

# The Curious Price of Distributional Robustness in Reinforcement Learning with a Generative Model

Laixi Shi\*    Gen Li†    Yuting Wei†    Yuxin Chen†    Matthieu Geist‡    Yuejie Chi\*  
CMU    UPenn    UPenn    UPenn    Google    CMU

May 25, 2023

## Abstract

This paper investigates model robustness in reinforcement learning (RL) to reduce the sim-to-real gap in practice. We adopt the framework of distributionally robust Markov decision processes (RMDPs), aimed at learning a policy that optimizes the worst-case performance when the deployed environment falls within a prescribed uncertainty set around the nominal MDP. Despite recent efforts, the sample complexity of RMDPs remained mostly unsettled regardless of the uncertainty set in use. It was unclear if distributional robustness bears any statistical consequences when benchmarked against standard RL.

Assuming access to a generative model that draws samples based on the nominal MDP, we characterize the sample complexity of RMDPs when the uncertainty set is specified via either the total variation (TV) distance or  $\chi^2$  divergence. The algorithm studied here is a model-based method called *distributionally robust value iteration*, which is shown to be near-optimal for the full range of uncertainty levels. Somewhat surprisingly, our results uncover that RMDPs are not necessarily easier or harder to learn than standard MDPs. The statistical consequence incurred by the robustness requirement depends heavily on the size and shape of the uncertainty set: in the case w.r.t. the TV distance, the minimax sample complexity of RMDPs is always smaller than that of standard MDPs; in the case w.r.t. the  $\chi^2$  divergence, the sample complexity of RMDPs can often far exceed the standard MDP counterpart.

**Keywords:** distributionally robust RL, robust Markov decision processes, sample complexity, distributionally robust value iteration, model-based RL

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Prior art and open questions	3
1.2	Main contributions	4
<b>2</b>	<b>Problem formulation</b>	<b>6</b>
<b>3</b>	<b>Model-based algorithm: distributionally robust value iteration</b>	<b>8</b>
<b>4</b>	<b>Theoretical guarantees: sample complexity analyses</b>	<b>9</b>
4.1	The case of TV distance: RMDPs are easier to learn than standard MDPs	9
4.2	The case of $\chi^2$ divergence: RMDPs can be harder than standard MDPs	11
<b>5</b>	<b>Other related works</b>	<b>13</b>
<b>6</b>	<b>Discussions</b>	<b>14</b>

---

\*Department of Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh, PA 15213, USA.

†Department of Statistics and Data Science, Wharton School, University of Pennsylvania, Philadelphia, PA 19104, USA.

‡Google Research, Brain Team.

<b>A Preliminaries</b>	<b>18</b>
A.1 Basic facts . . . . .	20
A.2 Properties of the robust Bellman operator . . . . .	20
A.3 Additional facts of the empirical robust MDP . . . . .	22
<b>B Proof of the upper bound with TV distance: Theorem 1</b>	<b>24</b>
B.1 Technical lemmas . . . . .	24
B.2 Proof of Theorem 1 . . . . .	24
B.3 Proof of the auxiliary lemmas . . . . .	33
<b>C Proof of the lower bound with TV distance: Theorem 2</b>	<b>42</b>
C.1 Construction of the hard problem instances . . . . .	42
C.2 Establishing the minimax lower bound . . . . .	44
C.3 Proof of the auxiliary facts . . . . .	45
<b>D Proof of the upper bound with <math>\chi^2</math> divergence: Theorem 3</b>	<b>49</b>
D.1 Proof of Theorem 3 . . . . .	49
D.2 Proof of the auxiliary lemmas . . . . .	51
<b>E Proof of the lower bound with <math>\chi^2</math> divergence: Theorem 4</b>	<b>55</b>
E.1 Construction of the hard problem instances . . . . .	55
E.2 Establishing the minimax lower bound . . . . .	56
E.3 Proof of the auxiliary facts . . . . .	58

# 1 Introduction

Reinforcement learning (RL) strives to learn desirable sequential decisions based on trial-and-error interactions with an unknown environment. As a fast-growing subfield of artificial intelligence, it has achieved remarkable success in a variety of domains such as large language model alignment (OpenAI, 2023; Ziegler et al., 2019), healthcare (Fatemi et al., 2021; Liu et al., 2019), and robotics and control (Kober et al., 2013; Mnih et al., 2013). Due to the unprecedented dimensionality of the state-action space, the issue of data efficiency inevitably lies at the core of modern RL practice. A large portion of recent efforts in RL has been directed towards designing sample-efficient algorithms and understanding the fundamental statistical bottleneck for a diverse range of RL scenarios.

While standard RL has been heavily investigated recently, its use can be significantly hampered in practice due to the sim-to-real gap; for instance, a policy learned in an ideal, nominal environment might fail catastrophically when the deployed environment is subject to small changes in task objectives or adversarial perturbations (Klopp et al., 2017; Mahmood et al., 2018; Zhang et al., 2020a). Consequently, in addition to maximizing the long-term cumulative reward, robustness emerges as another critical goal for RL, especially in high-stakes applications such as robotics, autonomous driving, clinical trials, financial investments, and so on. Towards achieving this, distributionally robust RL (Iyengar, 2005; Nilim and El Ghaoui, 2005), which leverages insights from distributionally robust optimization and supervised learning (Bertsimas et al., 2018; Blanchet and Murthy, 2019; Duchi and Namkoong, 2018; Gao, 2020; Rahimian and Mehrotra, 2019), becomes a natural yet versatile framework; the aim is to learn a policy that performs well even when the deployed environment deviates from the nominal one in the face of environment uncertainty.

In this paper, we pursue fundamental understanding about whether, and how, the choice of distributional robustness bears statistical implications in learning a desirable policy, through the lens of sample complexity. More concretely, imagine that one has access to a generative model (also called a simulator) that draws samples from a Markov decision processes (MDP) with a nominal transition kernel (Kearns and Singh, 1999). Standard RL aims to learn the optimal policy tailored to the nominal kernel, for which the minimax sample complexity limit has been fully settled (Azar et al., 2013b; Li et al., 2023b). In contrast, distributionally robust RL seeks to learn a more *robust* policy using the same set of samples, with the aim of optimizing the worst-case performance when the transition kernel is arbitrarily chosen from some *prescribed* uncertainty set

around the nominal kernel; this setting is frequently referred to as robust MDPs (RMDPs).<sup>1</sup> Clearly, the RMDP framework helps ensure that the performance of the learned policy does not fail catastrophically as long as the sim-to-real gap is not overly large. It is then natural to wonder how the robustness consideration impacts data efficiency: is there a statistical premium that one needs to pay in quest of additional robustness?

Compared with standard MDPs, the class of RMDPs encapsulates richer models, given that one is allowed to prescribe the shape and size of the uncertainty set. Oftentimes, the uncertainty set is hand-picked as a small ball surrounding the nominal kernel, with the size and shape of the ball specified by some distance-like metric  $\rho$  between probability distributions and some uncertainty level  $\sigma$ . To ensure tractability of solving RMDPs, the uncertainty set is often selected to obey certain structures. For instance, a number of prior works assumed that the uncertainty set can be decomposed as a product of independent uncertainty subsets over each state or state-action pair (Wiesemann et al., 2013; Zhou et al., 2021), dubbed as the  $s$ - and  $(s, a)$ -rectangularity, respectively. The current paper adopts the second choice by assuming  $(s, a)$ -rectangularity for the uncertainty set. An additional challenge with RMDPs arises from distribution shift, where the transition kernel drawn from the uncertainty set can be different from the nominal kernel. This challenge leads to complicated nonlinearity and nested optimization in the problem structure not present in standard MDPs.

## 1.1 Prior art and open questions

In this paper, we focus attention on RMDPs in the context of  $\gamma$ -discounted infinite-horizon setting, assuming access to a generative model. The uncertainty set considered herein is specified using one of the  $f$ -divergence metrics: the total variation (TV) distance and the  $\chi^2$  divergence. These two choices are motivated by their practical appeals: easy to implement, and already adopted by empirical RL (Lee et al., 2021).

A popular learning approach is model-based, which first estimates the nominal transition kernel using a plug-in estimator based on the collected samples, and then runs a planning algorithm (e.g., a robust variant of value iteration) on top of the estimated kernel. Despite the surge of recent activities, however, existing statistical guarantees for the above paradigm remained highly inadequate, as we shall elaborate on momentarily (see Table 1 and Table 2 respectively for a summary of existing results). For concreteness, let  $S$  be the size of the state space,  $A$  the size of the action space,  $\gamma$  the discount factor (so that the effective horizon is  $\frac{1}{1-\gamma}$ ), and  $\sigma$  the uncertainty level. We are interested in how the sample complexity — the number of samples needed for an algorithm to output a policy whose robust value function (the worst-case value over all the transition kernels in the uncertainty set) is at most  $\varepsilon$  away from the optimal robust one — scales with all these salient problem parameters.

- *Large gaps between existing upper and lower bounds.* There remained large gaps between the sample complexity upper and lower bounds established in prior literature, regardless of the divergence metric in use. For the case w.r.t. the TV distance, while the state-of-the-art upper bound (Clavier et al., 2023) and lower bound (Yang et al., 2022) coincide when the uncertainty level  $\sigma \lesssim 1 - \gamma$  is small,<sup>2</sup> the upper bound can be a factor of  $\frac{1}{(1-\gamma)^5}$  larger than the lower bound when  $\sigma$  approaches 1. The situation is even worse when it comes to the case w.r.t. the  $\chi^2$  divergence. More specifically, the state-of-the-art upper bound (Panaganti and Kalathil, 2022) scales quadratically with the size  $S$  of the state space and linearly with the uncertainty level  $\sigma$  when  $\sigma \gtrsim 1$ , while the lower bound (Yang et al., 2022) exhibits only linear scaling with  $S$  and is, in the meantime, inversely proportional to  $\sigma$ ; these lead to unbounded gaps between the upper and lower bounds as  $\sigma$  grows. *Can we hope to close these gaps for RMDPs?*
- *Benchmarking with standard MDPs.* Perhaps a more pressing issue is that, past works failed to provide an affirmative answer regarding how to benchmark the sample complexity of RMDPs with that of standard MDPs over the full range of uncertainty levels, given the large unresolved gaps mentioned above. In fact, prior works only achieved limited success in this regard — namely, demonstrating that the sample complexity for RMDPs is the same as that of standard MDPs in the case of TV distance

<sup>1</sup>While it is straightforward to incorporate additional uncertainty of the reward in our framework, we do not consider it here for simplicity, since the key challenge is to deal with the uncertainty of the transition kernel.

<sup>2</sup>Let  $\mathcal{X} := (S, A, \frac{1}{1-\gamma}, \sigma, \frac{1}{\varepsilon}, \frac{1}{\delta})$ . The notation  $f(\mathcal{X}) = O(g(\mathcal{X}))$  or  $f(\mathcal{X}) \lesssim g(\mathcal{X})$  indicates that there exists a universal constant  $C_1 > 0$  such that  $f \leq C_1 g$ , the notation  $f(\mathcal{X}) \gtrsim g(\mathcal{X})$  indicates that  $g(\mathcal{X}) = O(f(\mathcal{X}))$ , and the notation  $f(\mathcal{X}) \asymp g(\mathcal{X})$  indicates that  $f(\mathcal{X}) \lesssim g(\mathcal{X})$  and  $f(\mathcal{X}) \gtrsim g(\mathcal{X})$  hold simultaneously. Additionally, the notation  $\tilde{O}(\cdot)$  is defined in the same way as  $O(\cdot)$  except that it hides logarithmic factors.

Result type	Reference	Sample complexity	
		$0 < \sigma \lesssim 1 - \gamma$	$1 - \gamma \lesssim \sigma < 1$
Upper bound	Yang et al. (2022)	$\frac{S^2 A}{\sigma^2 (1-\gamma)^4 \varepsilon^2}$	
	Panaganti and Kalathil (2022)	$\frac{S^2 A}{(1-\gamma)^4 \varepsilon^2}$	
	Clavier et al. (2023)	$\frac{SA}{(1-\gamma)^3 \varepsilon^2}$	$\frac{SA}{(1-\gamma)^4 \varepsilon^2}$
	<b>This paper</b>	$\frac{SA}{(1-\gamma)^3 \varepsilon^2}$	$\frac{SA}{(1-\gamma)^2 \sigma \varepsilon^2}$
Lower bound	Yang et al. (2022)	$\frac{SA}{(1-\gamma)^3 \varepsilon^2}$	$\frac{SA(1-\gamma)}{\sigma^4 \varepsilon^2}$
	<b>This paper</b>	$\frac{SA}{(1-\gamma)^3 \varepsilon^2}$	$\frac{SA}{(1-\gamma)^2 \sigma \varepsilon^2}$

Table 1: Comparisons between our results and prior arts for finding an  $\varepsilon$ -optimal robust policy in the infinite-horizon RMDPs, with the uncertainty set measured w.r.t. the TV distance. Here,  $S$ ,  $A$ ,  $\gamma$ , and  $\sigma \in (0, 1)$  are the state space size, the action space size, the discount factor, and the uncertainty level, respectively, and all logarithmic factors are omitted in the table. Our results provide the first matching upper and lower bounds (up to log factors), improving upon all prior results.

when the uncertainty level satisfies  $\sigma \lesssim 1 - \gamma$ . For all the remaining situations, however, existing sample complexity upper (resp. lower) bounds are all larger (resp. smaller) than the sample size requirement for standard MDPs. As a consequence, it remains mostly unclear *whether learning RMDPs is harder or easier than learning standard MDPs*.

## 1.2 Main contributions

To address the aforementioned questions, this paper develops strengthened sample complexity upper bounds on learning RMDPs with the TV distance and  $\chi^2$  divergence in the infinite-horizon setting, using a model-based approach called distributionally robust value iteration (DRVI). Improved minimax lower bounds are also developed to help gauge the tightness of our upper bounds and enable benchmarking with standard MDPs. The novel analysis framework developed herein leads to new insights into the interplay between the geometry of uncertainty sets and statistical hardness.

**Sample complexity of RMDPs under the TV distance.** We summarize our results and compare them with past works in Table 1; see Figure 1(a) for a graphical illustration.

- **Minimax-optimal sample complexity.** We prove that DRVI reaches  $\varepsilon$  accuracy as soon as the sample complexity is on the order of

$$\tilde{O}\left(\frac{SA}{(1-\gamma)^2 \varepsilon^2} \min\left\{\frac{1}{1-\gamma}, \frac{1}{\sigma}\right\}\right)$$

for all  $\sigma \in (0, 1)$ , assuming that  $\varepsilon$  is small enough. In addition, a matching minimax lower bound (modulo some logarithmic factor) is established to guarantee the tightness of the upper bound over the full range of the uncertainty level. To the best of our knowledge, this is the *first* minimax-optimal sample complexity for RMDPs, which was previously unavailable regardless of the divergence metric in use.

- **RMDPs are easier to learn than standard MDPs under the TV distance.** Given the sample complexity  $\tilde{O}\left(\frac{SA}{(1-\gamma)^3 \varepsilon^2}\right)$  of standard MDPs (Li et al., 2023b), it can be seen that learning RMDPs

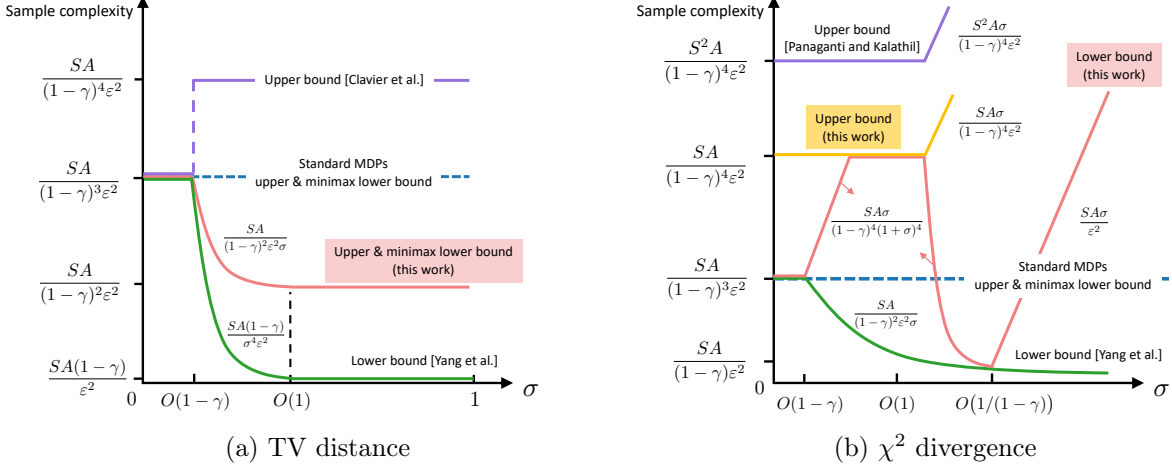


Figure 1: Illustrations of the obtained sample complexity upper and lower bounds for learning RMDPs with comparisons to state-of-the-art and the sample complexity of standard MDPs, where the uncertainty set is specified using the TV distance (a) and the  $\chi^2$  divergence (b).

under the TV distance is never harder than learning standard MDPs; more concretely, the sample complexity for RMDPs matches that of standard MDPs when  $\sigma \lesssim 1 - \gamma$ , and becomes smaller by a factor of  $\sigma/(1 - \gamma)$  when  $1 - \gamma \lesssim \sigma < 1$ . Therefore, in this case, distributional robustness comes almost for free, given that we do not need to collect more samples.

**Sample complexity of RMDPs under the  $\chi^2$  divergence.** We summarize our results and provide comparisons with prior works in Table 2; see Figure 1(b) for an illustration.

- **Near-optimal sample complexity.** We demonstrate that DRVI yields  $\varepsilon$  accuracy as soon as the sample complexity is on the order of

$$\tilde{O}\left(\frac{SA(1 + \sigma)}{(1 - \gamma)^4 \varepsilon^2}\right)$$

for all  $\sigma \in (0, \infty)$ , which is the first sample complexity in this setting that scales linearly in the size  $S$  of the state space; in other words, our theory breaks the quadratic scaling bottleneck that was present in prior works (Panaganti and Kalathil, 2022; Yang et al., 2022). We have also developed a strengthened lower bound that is optimized by leveraging the geometry of the uncertainty set under different ranges of  $\sigma$ . Our theory is tight when  $\sigma \asymp 1$ , and is otherwise loose by at most a polynomial factor of the effective horizon  $1/(1 - \gamma)$  (regardless of the uncertainty level  $\sigma$ ). This significantly improves upon prior results (as there exists an unbounded gap between prior upper and lower bounds as  $\sigma \rightarrow \infty$ ).

- **RMDPs can be harder to learn than standard MDPs under the  $\chi^2$  divergence.** Somewhat surprisingly, our improved lower bound suggests that RMDPs in this case can be much harder to learn than standard MDPs, at least for a certain range of uncertainty levels. We single out two regimes of particular interest. Firstly, when  $\sigma \asymp 1$ , the sample size requirement of RMDPs is on the order of  $\frac{SA}{(1 - \gamma)^4 \varepsilon^2}$  (up to log factor), which is provably larger than the one for standard MDPs by a factor of  $\frac{1}{1 - \gamma}$ . Secondly, the lower bound continues to increase as  $\sigma$  grows and exceeds the sample complexity of standard MDPs when  $\sigma \gtrsim \frac{1}{(1 - \gamma)^3}$ .

In sum, our sample complexity bounds not only strengthen the prior art in the development of both upper and lower bounds, but also unveil that the additional robustness consideration might affect the sample complexity in a somewhat surprising manner. As it turns out, RMDPs are not necessarily harder nor easier to learn than standard MDPs; the conclusion is far more nuanced and highly dependent on both the size and shape of the uncertainty set. This constitutes a curious phenomenon that has not been elucidated in prior analyses.

Result type	Reference	Sample complexity		
		$0 < \sigma \lesssim 1 - \gamma$	$1 - \gamma \lesssim \sigma \lesssim \frac{1}{1-\gamma}$	$\sigma \gtrsim \frac{1}{1-\gamma}$
Upper bound	Panaganti and Kalathil (2022)	$\frac{S^2 A(1+\sigma)}{(1-\gamma)^4 \varepsilon^2}$		
	Yang et al. (2022)	$\frac{S^2 A(1+\sigma)^2}{(\sqrt{1+\sigma}-1)^2 (1-\gamma)^4 \varepsilon^2}$		
	<b>This paper</b>	$\frac{SA(1+\sigma)}{(1-\gamma)^4 \varepsilon^2}$		
Lower bound	Yang et al. (2022)	$\frac{SA}{(1-\gamma)^3 \varepsilon^2}$	$\frac{SA}{(1-\gamma)^2 \sigma \varepsilon^2}$	
	<b>This paper</b>	$\frac{SA}{(1-\gamma)^3 \varepsilon^2}$	$\frac{SA\sigma}{(1-\gamma)^4 (1+\sigma)^4 \varepsilon^2}$	$\frac{SA\sigma}{\varepsilon^2}$

Table 2: Comparisons between our results and prior art on finding an  $\varepsilon$ -optimal robust policy in the infinite-horizon RMDPs, with the uncertainty set measured w.r.t. the  $\chi^2$  divergence. Here,  $S$ ,  $A$ ,  $\gamma$ , and  $\sigma \in (0, \infty)$  are the state space size, the action space size, the discount factor, and the uncertainty level, respectively, and all logarithmic factors are omitted in the table. Improving upon all prior results, our theory is tight (up to log factors) when  $\sigma \asymp 1$ , and otherwise loose by no more than a polynomial factor in  $1/(1-\gamma)$ .

**Notation and paper organization.** Throughout this paper, we denote by  $\Delta(\mathcal{S})$  the probability simplex over a set  $\mathcal{S}$  and  $x = [x(s, a)]_{(s,a) \in \mathcal{S} \times \mathcal{A}} \in \mathbb{R}^{SA}$  (resp.  $x = [x(s)]_{s \in \mathcal{S}} \in \mathbb{R}^S$ ) as any vector that constitutes certain values for each state-action pair (resp. state). In addition, we denote by  $x \circ y = [x(s) \cdot y(s)]_{s \in \mathcal{S}}$  the Hadamard product of any two vectors  $x, y \in \mathbb{R}^S$ .

The remainder of this paper is structured as follows. Section 2 presents the background about discounted infinite-horizon standard MDPs and formulates distributionally robust MDPs. In Section 3, a model-based approach is introduced, tailored to both the TV distance and the  $\chi^2$  divergence. Both upper and lower bounds on the sample complexity are developed in Section 4, covering both divergence metrics. We then summarize several additional related works in Section 5 and conclude the main paper with further discussions in Section 6. The proof details are deferred to the appendix.

## 2 Problem formulation

In this section, we formulate distributionally robust Markov decision processes (RMDPs) in the discounted infinite-horizon setting, introduce the sampling mechanism, and describe our goal.

**Standard MDPs.** To begin, we first introduce the standard Markov decision processes (MDPs), which facilitate the understanding of RMDPs. A discounted infinite-horizon MDP is represented by  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \gamma, P, r)$ , where  $\mathcal{S} = \{1, \dots, S\}$  and  $\mathcal{A} = \{1, \dots, A\}$  are the finite state and action spaces, respectively,  $\gamma \in [0, 1]$  is the discounted factor,  $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$  denotes the probability transition kernel, and  $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  is the immediate reward function which is assumed to be deterministic. A policy is denoted by  $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ , which specifies the action selection probability over the action space in any state. When the policy is deterministic, we overload the notation and refer to  $\pi(s)$  as the action selected by policy  $\pi$  in state  $s$ . To characterize the cumulative reward, the value function  $V^{\pi, P}$  for any policy  $\pi$  under the transition kernel  $P$  is defined by

$$\forall s \in \mathcal{S} : \quad V^{\pi, P}(s) := \mathbb{E}_{\pi, P} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s \right], \quad (1)$$

where the expectation is taken over the randomness of the trajectory  $\{s_t, a_t\}_{t=0}^{\infty}$  generated by executing policy  $\pi$  under the transition kernel  $P$ , namely,  $a_t \sim \pi(\cdot \mid s_t)$  and  $s_{t+1} \sim P(\cdot \mid s_t, a_t)$  for all  $t \geq 0$ . Similarly,

the Q-function  $Q^{\pi,P}$  associated with any policy  $\pi$  under the transition kernel  $P$  is defined as

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad Q^{\pi,P}(s, a) := \mathbb{E}_{\pi,P} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, a_0 = a \right], \quad (2)$$

where the expectation is again taken over the randomness of the trajectory under policy  $\pi$ .

**Distributionally robust MDPs.** We now introduce the distributionally robust MDP (RMDP) tailored to the discounted infinite-horizon setting, denoted by  $\mathcal{M}_{\text{rob}} = \{\mathcal{S}, \mathcal{A}, \gamma, \mathcal{U}_{\rho}^{\sigma}(P^0), r\}$ , where  $\mathcal{S}, \mathcal{A}, \gamma, r$  are identical to those in the standard MDP. A key distinction from the standard MDP is that: rather than assuming a fixed transition kernel  $P$ , it allows the transition kernel to be chosen arbitrarily from a prescribed uncertainty set  $\mathcal{U}_{\rho}^{\sigma}(P^0)$  centered around a *nominal* kernel  $P^0 : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ , where the uncertainty set is specified using some distance metric  $\rho$  of radius  $\sigma > 0$ . In particular, given the nominal transition kernel  $P^0$  and some uncertainty level  $\sigma$ , the uncertainty set—with the divergence metric  $\rho : \Delta(\mathcal{S}) \times \Delta(\mathcal{S}) \rightarrow \mathbb{R}^+$ —is specified as

$$\mathcal{U}_{\rho}^{\sigma}(P^0) := \otimes \mathcal{U}_{\rho}^{\sigma}(P_{s,a}^0) \quad \text{with} \quad \mathcal{U}_{\rho}^{\sigma}(P_{s,a}^0) := \{P_{s,a} \in \Delta(\mathcal{S}) : \rho(P_{s,a}, P_{s,a}^0) \leq \sigma\}, \quad (3)$$

where we denote a vector of the transition kernel  $P$  or  $P^0$  at state-action pair  $(s, a)$  respectively as

$$P_{s,a} := P(\cdot \mid s, a) \in \mathbb{R}^{1 \times \mathcal{S}}, \quad P_{s,a}^0 := P^0(\cdot \mid s, a) \in \mathbb{R}^{1 \times \mathcal{S}}. \quad (4)$$

In other words, the uncertainty is imposed in a decoupled manner for each state-action pair, obeying the so-called  $(s, a)$ -rectangularity (Wiesemann et al., 2013; Zhou et al., 2021).

In RMDPs, we are interested in the worst-case performance of a policy  $\pi$  over all the possible transition kernels in the uncertainty set. This is measured by the *robust value function*  $V^{\pi,\sigma}$  and the *robust Q-function*  $Q^{\pi,\sigma}$  in  $\mathcal{M}_{\text{rob}}$ , defined respectively as

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad V^{\pi,\sigma}(s) := \inf_{P \in \mathcal{U}_{\rho}^{\sigma}(P^0)} V^{\pi,P}(s), \quad Q^{\pi,\sigma}(s, a) := \inf_{P \in \mathcal{U}_{\rho}^{\sigma}(P^0)} Q^{\pi,P}(s, a). \quad (5)$$

**Optimal robust policy and robust Bellman operator.** As a generalization of properties of standard MDPs, it is well-known that there exists at least one deterministic policy that maximizes the robust value function (resp. robust Q-function) simultaneously for all states (resp. state-action pairs) (Iyengar, 2005; Nilim and El Ghaoui, 2005). Therefore, we denote the *optimal robust value function* (resp. *optimal robust Q-function*) as  $V^{*,\sigma}$  (resp.  $Q^{*,\sigma}$ ), and the optimal robust policy as  $\pi^*$ , which satisfy

$$\forall s \in \mathcal{S}: \quad V^{*,\sigma}(s) := V^{\pi^*,\sigma}(s) = \max_{\pi} V^{\pi,\sigma}(s), \quad (6a)$$

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad Q^{*,\sigma}(s, a) := Q^{\pi^*,\sigma}(s, a) = \max_{\pi} Q^{\pi,\sigma}(s, a). \quad (6b)$$

A key machinery in RMDPs is a generalization of Bellman’s optimality principle, encapsulated in the following *robust Bellman consistency equation* (resp. *robust Bellman optimality equation*):

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad Q^{\pi,\sigma}(s, a) = r(s, a) + \gamma \inf_{P \in \mathcal{U}_{\rho}^{\sigma}(P_{s,a}^0)} \mathcal{P}V^{\pi,\sigma}, \quad (7a)$$

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad Q^{*,\sigma}(s, a) = r(s, a) + \gamma \inf_{P \in \mathcal{U}_{\rho}^{\sigma}(P_{s,a}^0)} \mathcal{P}V^{*,\sigma}. \quad (7b)$$

The robust Bellman operator (Iyengar, 2005; Nilim and El Ghaoui, 2005) is denoted by  $\mathcal{T}^{\sigma}(\cdot) : \mathbb{R}^{\mathcal{S}\mathcal{A}} \rightarrow \mathbb{R}^{\mathcal{S}\mathcal{A}}$  and defined as follows:

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad \mathcal{T}^{\sigma}(Q)(s, a) := r(s, a) + \gamma \inf_{P \in \mathcal{U}_{\rho}^{\sigma}(P_{s,a}^0)} \mathcal{P}V, \quad \text{with} \quad V(s) := \max_a Q(s, a). \quad (8)$$

Given that  $Q^{*,\sigma}$  is the unique fixed point of  $\mathcal{T}^{\sigma}$ , one can recover the optimal robust value function and Q-function using a procedure termed *distributionally robust value iteration* (DRVI). Generalizing the standard value iteration, DRVI starts from some given initialization and recursively applies the robust Bellman operator until convergence. As has been shown previously, this procedure converges rapidly due to the  $\gamma$ -contraction property of  $\mathcal{T}^{\sigma}$  w.r.t. the  $\ell_{\infty}$  norm (Iyengar, 2005; Nilim and El Ghaoui, 2005).



**Specification of the divergence  $\rho$ .** We consider two popular choices of the uncertainty set measured in terms of two different  $f$ -divergence metric: the total variation distance and the  $\chi^2$  divergence, given respectively by (Tsybakov, 2009)

$$\rho_{\text{TV}}(P_{s,a}, P_{s,a}^0) := \frac{1}{2} \|P_{s,a} - P_{s,a}^0\|_1 = \frac{1}{2} \sum_{s' \in \mathcal{S}} P^0(s' | s, a) \left| 1 - \frac{P(s' | s, a)}{P^0(s' | s, a)} \right|, \quad (9)$$

$$\rho_{\chi^2}(P_{s,a}, P_{s,a}^0) := \sum_{s' \in \mathcal{S}} P^0(s' | s, a) \left( 1 - \frac{P(s' | s, a)}{P^0(s' | s, a)} \right)^2. \quad (10)$$

Note that  $\rho_{\text{TV}}(P_{s,a}, P_{s,a}^0) \in [0, 1]$  and  $\rho_{\chi^2}(P_{s,a}, P_{s,a}^0) \in [0, \infty)$  in general. As we shall see shortly, these two choices of divergence metrics result in drastically different messages when it comes to sample complexities.

**Sampling mechanism: a generative model.** Following Panaganti and Kalathil (2022); Zhou et al. (2021), we assume access to a generative model or a simulator (Kearns and Singh, 1999), which allows us to collect  $N$  independent samples for each state-action pair generated based on the *nominal* kernel  $P^0$ :

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}, \quad s_{i,s,a} \stackrel{i.i.d}{\sim} P^0(\cdot | s, a), \quad i = 1, 2, \dots, N. \quad (11)$$

The total sample size is, therefore,  $NSA$ .

**Goal.** Given the collected samples, the task is to learn the robust optimal policy for the RMDP — w.r.t. some prescribed uncertainty set  $\mathcal{U}^\sigma(P^0)$  around the nominal kernel — using as few samples as possible. Specifically, given some target accuracy level  $\varepsilon > 0$ , the goal is to seek an  $\varepsilon$ -optimal robust policy  $\hat{\pi}$  obeying

$$\forall s \in \mathcal{S}: \quad V^{\star, \sigma}(s) - V^{\hat{\pi}, \sigma}(s) \leq \varepsilon. \quad (12)$$

### 3 Model-based algorithm: distributionally robust value iteration

We consider a model-based approach tailored to RMDPs, which first constructs an empirical nominal transition kernel based on the collected samples, and then applies distributionally robust value iteration (DRVI) to compute an optimal robust policy.

**Empirical nominal kernel.** The empirical nominal transition kernel  $\hat{P}^0 \in \mathbb{R}^{\mathcal{S} \times \mathcal{S}}$  can be constructed on the basis of the empirical frequency of state transitions, i.e.,

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad \hat{P}^0(s' | s, a) := \frac{1}{N} \sum_{i=1}^N \mathbb{1}\{s_{i,s,a} = s'\}, \quad (13)$$

which leads to an empirical RMDP  $\widehat{\mathcal{M}}_{\text{rob}} = \{\mathcal{S}, \mathcal{A}, \gamma, \mathcal{U}_\rho^\sigma(\hat{P}^0), r\}$ . Analogously, we can define the corresponding robust value function (resp. robust Q-function) of policy  $\pi$  in  $\widehat{\mathcal{M}}_{\text{rob}}$  as  $\widehat{V}^{\pi, \sigma}$  (resp.  $\widehat{Q}^{\pi, \sigma}$ ) (cf. (6)). In addition, we denote the corresponding *optimal robust policy* as  $\hat{\pi}^\star$  and the *optimal robust value function* (resp. *optimal robust Q-function*) as  $\widehat{V}^{\star, \sigma}$  (resp.  $\widehat{Q}^{\star, \sigma}$ ) (cf. (7)), which satisfies the robust Bellman optimality equation:

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad \widehat{Q}^{\star, \sigma}(s, a) = r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}_\rho^\sigma(\hat{P}_{s,a}^0)} \mathcal{P} \widehat{V}^{\star, \sigma}. \quad (14)$$

Equipped with  $\hat{P}^0$ , we can define the empirical robust Bellman operator  $\widehat{\mathcal{T}}^\sigma$  as

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad \widehat{\mathcal{T}}^\sigma(Q)(s, a) := r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}_\rho^\sigma(\hat{P}_{s,a}^0)} \mathcal{P} V, \quad \text{with} \quad V(s) := \max_a Q(s, a). \quad (15)$$



---

**Algorithm 1:** Distributionally robust value iteration (DRVI) for infinite-horizon RMDPs.

---

```

1 input: empirical nominal transition kernel  $\hat{P}^0$ ; reward function  $r$ ; uncertainty level  $\sigma$ ; number of
   iterations  $T$ .
2 initialization:  $\hat{Q}_0(s, a) = 0$ ,  $\hat{V}_0(s) = 0$  for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ .
3 for  $t = 1, 2, \dots, T$  do
4   for  $s \in \mathcal{S}, a \in \mathcal{A}$  do
5      $\hat{Q}_t(s, a)$  according to (16);
6   for  $s \in \mathcal{S}$  do
7      $\hat{V}_t(s) = \max_a \hat{Q}_t(s, a)$ ;
8 output:  $\hat{Q}_T$ ,  $\hat{V}_T$  and  $\hat{\pi}$  obeying  $\hat{\pi}(s) := \arg \max_a \hat{Q}_T(s, a)$ .
```

---

**DRVI: distributionally robust value iteration.** To compute the fixed point of  $\hat{\mathcal{T}}^\sigma$ , we introduce distributionally robust value iteration (DRVI), which is summarized in Algorithm 1. From an initialization  $\hat{Q}_0 = 0$ , the update rule at the  $t$ -th ( $t \geq 1$ ) iteration can be formulated as:

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad \hat{Q}_t(s, a) = \hat{\mathcal{T}}^\sigma(\hat{Q}_{t-1})(s, a) = r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}_\rho^\sigma(\hat{P}_{s,a}^0)} \mathcal{P}\hat{V}_{t-1}, \quad (16)$$

where  $\hat{V}_{t-1}(s) = \max_a \hat{Q}_{t-1}(s, a)$  for all  $s \in \mathcal{S}$ . However, directly solving (16) is computationally expensive since it involves optimization over an  $S$ -dimensional probability simplex at each iteration, especially when the dimension of the state space  $\mathcal{S}$  is large. Fortunately, in view of strong duality (Iyengar, 2005), (16) can be equivalently solved using its dual problem, which concerns optimizing a *scalar* dual variable and thus can be solved efficiently. The specific form of the dual problem depends on the choice of the divergence  $\rho$ , which we shall discuss separately in Appendix A.2. To complete the description, we output the greedy policy of the final  $Q$ -estimate  $\hat{Q}_T$  as the final policy  $\hat{\pi}$ , namely,

$$\forall s \in \mathcal{S}: \quad \hat{\pi}(s) = \arg \max_a \hat{Q}_T(s, a). \quad (17)$$

Encouragingly, the iterates  $\{\hat{Q}_t\}_{t \geq 0}$  of DRVI converge linearly to the fixed point  $\hat{Q}^{*,\sigma}$ , owing to the appealing  $\gamma$ -contraction property of  $\hat{\mathcal{T}}^\sigma$ .

## 4 Theoretical guarantees: sample complexity analyses

We now present our main results, which concern the sample complexities of learning RMDPs when the uncertainty set is specified using the TV distance or the  $\chi^2$  divergence. Somewhat surprisingly, different choices of the uncertainty set can lead to dramatically different consequences in the sample size requirement.

### 4.1 The case of TV distance: RMDPs are easier to learn than standard MDPs

We start with the case where the uncertainty set is measured via the TV distance. The following theorem, whose proof is deferred to Appendix B, develops an upper bound on the sample complexity of DRVI in order to return an  $\varepsilon$ -optimal robust policy. The key challenge of the analysis lies in careful control of the robust value function  $V^{\pi,\sigma}$  as a function of the uncertainty level  $\sigma$ .

**Theorem 1** (Upper bound under TV distance). *Let the uncertainty set be  $\mathcal{U}_\rho^\sigma(\cdot) = \mathcal{U}_{\text{TV}}^\sigma(\cdot)$ , as specified by the TV distance (9). Consider any discount factor  $\gamma \in [\frac{1}{4}, 1)$ , uncertainty level  $\sigma \in (0, 1)$ , and  $\delta \in (0, 1)$ . Let  $\hat{\pi}$  be the output policy of Algorithm 1 after  $T = C_1 \log\left(\frac{N}{1-\gamma}\right)$  iterations. Then with probability at least  $1 - \delta$ , one has*

$$\forall s \in \mathcal{S}: \quad V^{*,\sigma}(s) - V^{\hat{\pi},\sigma}(s) \leq \varepsilon \quad (18)$$

for any  $\varepsilon \in \left(0, \sqrt{1/\max\{1-\gamma, \sigma\}}\right]$ , as long as the total number of samples obeys

$$NSA \geq \frac{C_2 SA}{(1-\gamma)^2 \max\{1-\gamma, \sigma\} \varepsilon^2} \log \left( \frac{SAN}{(1-\gamma)\delta} \right). \quad (19)$$

Here,  $C_1, C_2 > 0$  are some large enough universal constants.

**Remark 1.** Note that Theorem 1 is not only valid when invoking Algorithm 1. In fact, the theorem holds for any oracle planning algorithm (designed based on the empirical transitions  $\hat{P}^0$ ) whose output policy  $\hat{\pi}$  obeys

$$\|\hat{V}^{\star, \sigma} - \hat{V}^{\hat{\pi}, \sigma}\|_{\infty} \leq O \left( \frac{(1-\gamma)^2}{N} \log \left( \frac{SAN}{(1-\gamma)\delta} \right) \right). \quad (20)$$

Before discussing the implications of Theorem 1, we present a matching minimax lower bound that confirms the tightness and optimality of the upper bound, which in turn pins down the sample complexity requirement for learning RMDPs with TV distance. The proof is based on constructing new hard instances inspired by the asymmetric structure of RMDPs, with the details postponed to Appendix C.

**Theorem 2** (Lower bound under TV distance). *Consider any tuple  $(S, A, \gamma, \sigma, \varepsilon)$  obeying  $\sigma \in (0, 1-c_0]$  with  $0 < c_0 \leq \frac{1}{8}$  being any small enough positive constant,  $\gamma \in [\frac{1}{2}, 1)$ , and  $\varepsilon \in (0, \frac{c_0}{256(1-\gamma)})$ . We can construct two infinite-horizon RMDPs  $\mathcal{M}_0, \mathcal{M}_1$  defined by the uncertainty set  $\mathcal{U}_{\rho}^{\sigma}(\cdot) = \mathcal{U}_{TV}^{\sigma}(\cdot)$ , an initial state distribution  $\varphi$ , and a dataset with  $N$  independent samples for each state-action pair over the nominal transition kernel (for  $\mathcal{M}_0$  and  $\mathcal{M}_1$  respectively), such that*

$$\inf_{\hat{\pi}} \max \left\{ \mathbb{P}_0(V^{\star, \sigma}(\varphi) - V^{\hat{\pi}, \sigma}(\varphi) > \varepsilon), \mathbb{P}_1(V^{\star, \sigma}(\varphi) - V^{\hat{\pi}, \sigma}(\varphi) > \varepsilon) \right\} \geq \frac{1}{8},$$

provided that

$$NSA \leq \frac{c_0 SA \log 2}{8192(1-\gamma)^2 \max\{1-\gamma, \sigma\} \varepsilon^2}.$$

Here, the infimum is taken over all estimators  $\hat{\pi}$ , and  $\mathbb{P}_0$  (resp.  $\mathbb{P}_1$ ) denotes the probability when the RMDP is  $\mathcal{M}_0$  (resp.  $\mathcal{M}_1$ ).

Below, we interpret the above theorems and highlight several key implications about the sample complexity requirements for learning RMDPs for the case w.r.t. the TV distance.

**Near minimax-optimal sample complexity.** Theorem 1 shows that the total number of samples required for DRVI (or any oracle planning algorithm claimed in Remark 1) to yield  $\varepsilon$ -accuracy is

$$\tilde{O} \left( \frac{SA}{(1-\gamma)^2 \max\{1-\gamma, \sigma\} \varepsilon^2} \right). \quad (21)$$

Taken together with the minimax lower bound asserted by Theorem 2, this confirms the near optimality of the sample complexity (up to some logarithmic factor) almost over the full range of the uncertainty level  $\sigma$ . Importantly, this sample complexity scales linearly with the size of the state-action space, and is inversely proportional to  $\sigma$  in the regime where  $\sigma \gtrsim 1-\gamma$ .

**RMDPs is easier than standard MDPs with TV distance.** Recall that the sample complexity requirement for learning standard MDPs with a generative model is (Agarwal et al., 2020; Azar et al., 2013a; Li et al., 2023b)

$$\tilde{O} \left( \frac{SA}{(1-\gamma)^3 \varepsilon^2} \right) \quad (22)$$

in order to yield  $\varepsilon$  accuracy. Comparing this with the sample complexity requirement in (21) for RMDPs under the TV distance, we confirm that the latter is at least as easy as — if not easier than — standard MDPs. In particular, when  $\sigma \lesssim 1-\gamma$  is small, the sample complexity of RMDPs is the same as that of

standard MDPs as in (22), which is as anticipated since the RMDP reduces to the standard MDP when  $\sigma = 0$ . On the other hand, when  $1 - \gamma \lesssim \sigma < 1$ , the sample complexity of RMDPs simplifies to

$$\tilde{O}\left(\frac{SA}{(1-\gamma)^2\sigma\varepsilon^2}\right), \quad (23)$$

which is smaller than that of standard MDPs by a factor of  $\sigma/(1-\gamma)$ .

**Comparison with state-of-the-art bounds.** While the state-of-the-art sample complexity upper bound derived in Clavier et al. (2023) is tight when  $\sigma$  is small (i.e.,  $\sigma \lesssim 1-\gamma$ ), the sample complexity bound therein scales as  $\tilde{O}(\frac{SA}{(1-\gamma)^4\varepsilon^2})$  in the regime where  $1-\gamma \lesssim \sigma < 1$ . Consequently, this is worse than our result by a factor of

$$\frac{\sigma}{(1-\gamma)^2} \in \left(\frac{1}{1-\gamma}, \frac{1}{(1-\gamma)^2}\right).$$

Turning to the lower bound side, Yang et al. (2022) developed a lower bound for RMDPs under the TV distance, which scales as

$$\tilde{O}\left(\frac{SA(1-\gamma)}{\varepsilon^2} \min\left\{\frac{1}{(1-\gamma)^4}, \frac{1}{\sigma^4}\right\}\right).$$

Clearly, this is worse than ours by a factor of  $\frac{\sigma^3}{(1-\gamma)^3} \in (1, \frac{1}{(1-\gamma)^3})$  in the regime where  $1-\gamma \lesssim \sigma < 1$ .

## 4.2 The case of $\chi^2$ divergence: RMDPs can be harder than standard MDPs

We now switch attention to the case when the uncertainty set is measured via the  $\chi^2$  divergence. The theorem below presents an upper bound on the sample complexity for this case, whose proof is deferred to Appendix D.

**Theorem 3** (Upper bound under  $\chi^2$  divergence). *Let the uncertainty set be  $\mathcal{U}_\rho^\sigma(\cdot) = \mathcal{U}_{\chi^2}^\sigma(\cdot)$ , as specified using the  $\chi^2$  divergence (10). Consider any uncertainty level  $\sigma \in (0, \infty)$ ,  $\gamma \in [1/4, 1)$  and  $\delta \in (0, 1)$ . With probability at least  $1-\delta$ , the output policy  $\hat{\pi}$  from Algorithm 1 with at most  $T = c_1 \log(\frac{N}{1-\gamma})$  iterations yields*

$$\forall s \in \mathcal{S}: \quad V^{*,\sigma}(s) - V^{\hat{\pi},\sigma}(s) \leq \varepsilon \quad (24)$$

for any  $\varepsilon \in (0, \frac{1}{1-\gamma}]$ , as long as the total number of samples obeying

$$NSA \geq \frac{c_2 SA(1+\sigma)}{(1-\gamma)^4 \varepsilon^2} \log\left(\frac{SAN}{\delta}\right). \quad (25)$$

Here,  $c_1, c_2 > 0$  are some large enough universal constants.

**Remark 2.** Akin to Remark 1, the sample complexity derived in Theorem 3 continues to hold for any oracle planning algorithm that outputs a policy  $\hat{\pi}$  obeying  $\|\hat{V}^{*,\sigma} - \hat{V}^{\hat{\pi},\sigma}\|_\infty \leq O\left(\frac{\log(\frac{SAN}{(1-\gamma)\delta})}{N^2}\right)$ .

In addition, in order to gauge the tightness of Theorem 3 and understand the minimal sample complexity requirement under the  $\chi^2$  divergence, we further develop a minimax lower bound as follows; the proof is deferred to Appendix E.

**Theorem 4** (Lower bound under  $\chi^2$  divergence). *Consider any  $(S, A, \gamma, \sigma, \varepsilon)$  obeying  $\gamma \in [\frac{3}{4}, 1)$ ,  $\sigma \in (0, \infty)$ , and*

$$\varepsilon \leq c_3 \begin{cases} \frac{1}{1-\gamma} & \text{if } \sigma \in (0, \frac{1-\gamma}{4}) \\ \max\left\{\frac{1}{(1+\sigma)(1-\gamma)}, 1\right\} & \text{if } \sigma \in [\frac{1-\gamma}{4}, \infty) \end{cases} \quad (26)$$

for some small universal constant  $c_3 > 0$ . Then we can construct two infinite-horizon RMDPs  $\mathcal{M}_0, \mathcal{M}_1$  defined by the uncertainty set  $\mathcal{U}_\rho^\sigma(\cdot) = \mathcal{U}_{\chi^2}^\sigma(\cdot)$ , an initial state distribution  $\varphi$ , and a dataset with  $N$  independent samples per  $(s, a)$  pair over the nominal transition kernel (for  $\mathcal{M}_0$  and  $\mathcal{M}_1$  respectively), such that

$$\inf_{\hat{\pi}} \max \left\{ \mathbb{P}_0(V^{*,\sigma}(\varphi) - V^{\hat{\pi},\sigma}(\varphi) > \varepsilon), \mathbb{P}_1(V^{*,\sigma}(\varphi) - V^{\hat{\pi},\sigma}(\varphi) > \varepsilon) \right\} \geq \frac{1}{8}, \quad (27)$$

provided that the total number of samples

$$NSA \leq c_4 \begin{cases} \frac{SA}{(1-\gamma)^3 \varepsilon^2} & \text{if } \sigma \in (0, \frac{1-\gamma}{4}) \\ \frac{\sigma SA}{\min\{1, (1-\gamma)^4(1+\sigma)^4\} \varepsilon^2} & \text{if } \sigma \in [\frac{1-\gamma}{4}, \infty) \end{cases} \quad (28)$$

for some universal constant  $c_4 > 0$ .

We are now positioned to single out several key implications of the above theorems.

**Nearly tight sample complexity.** In order to achieve  $\varepsilon$ -accuracy for RMDPs under the  $\chi^2$  divergence, Theorem 3 asserts that a total number of samples on the order of

$$\tilde{O}\left(\frac{SA(1+\sigma)}{(1-\gamma)^4 \varepsilon^2}\right). \quad (29)$$

is sufficient for DRVI (or any other oracle planning algorithm as discussed in Remark 2). Taking this together with the minimax lower bound in Theorem 4 confirms that the sample complexity is near-optimal — up to a polynomial factor of the effective horizon  $\frac{1}{1-\gamma}$  — over the entire range of the uncertainty level  $\sigma$ . In particular,

- when  $\sigma \asymp 1$ , our sample complexity  $\tilde{O}\left(\frac{SA}{(1-\gamma)^4 \varepsilon^2}\right)$  is sharp and matches the minimax lower bound;
- when  $\sigma \gtrsim \frac{1}{(1-\gamma)^3}$ , our sample complexity correctly predicts the linear dependency with  $\sigma$ , suggesting that more samples are needed when one wishes to account for a larger  $\chi^2$ -based uncertainty sets.

**RMDPs can be much harder to learn than standard MDPs with  $\chi^2$  divergence.** The minimax lower bound developed in Theorem 4 exhibits a curious non-monotonic behavior of the sample size requirement over the entire range of the uncertainty level  $\sigma \in (0, \infty)$  when the uncertainty set is measured via the  $\chi^2$  divergence. When  $\sigma \lesssim 1 - \gamma$ , the lower bound reduces to

$$\tilde{O}\left(\frac{SA}{(1-\gamma)^3 \varepsilon^2}\right),$$

which matches with that of standard MDPs, as  $\sigma = 0$  corresponds to standard MDP. However, two additional regimes are worth calling out:

$$\begin{aligned} 1 - \gamma \lesssim \sigma \lesssim \frac{1}{(1-\gamma)^{1/3}} : \quad & \tilde{O}\left(\frac{SA}{(1-\gamma)^4 \varepsilon^2} \min\left\{\sigma, \frac{1}{\sigma^3}\right\}\right), \\ \sigma \gtrsim \frac{1}{(1-\gamma)^3} : \quad & \tilde{O}\left(\frac{SA\sigma}{\varepsilon^2}\right), \end{aligned}$$

both of which are *greater* than that of standard MDPs, indicating learning RMDPs under the  $\chi^2$  divergence can be much harder.

**Comparison with state-of-the-art bounds.** Our upper bound significantly improves over the prior art  $\tilde{O}\left(\frac{S^2 A(1+\sigma)}{(1-\gamma)^4 \varepsilon^2}\right)$  of Panaganti and Kalathil (2022) by a factor of  $S$ , and provides the *first* finite-sample complexity that scales *linearly* with respect to  $S$  for discounted infinite-horizon RMDPs, which typically exhibit more complicated statistical dependencies than the finite-horizon counterpart. On the other hand, Yang et al. (2022) established a lower bound on the order of  $\tilde{O}\left(\frac{SA}{(1-\gamma)^2 \sigma \varepsilon^2}\right)$  when  $\sigma \gtrsim 1 - \gamma$ , which is always smaller than the requirement of standard MDPs, and diminishes when  $\sigma$  grows. Consequently, Yang et al. (2022) does not lead to the rigorous justification that RMDPs can be much harder than standard MDPs, nor the correct linear scaling of the sample size as  $\sigma$  grows.

## 5 Other related works

This section briefly discusses a small sample of other related works. We limit our discussions primarily to provable RL algorithms in the tabular setting with finite state and action spaces, which are most related to the current paper.

**Finite-sample guarantees for standard RL.** A surge of recent research has utilized the toolkit from high-dimensional probability/statistics to investigate the performance of standard RL algorithms in non-asymptotic settings. There has been a considerable amount of research into non-asymptotic sample analysis of standard RL for a variety of settings; partial examples include, but are not limited to, the works via probably approximately correct (PAC) bounds for the generative model setting (Agarwal et al., 2020; Azar et al., 2013b; Beck and Srikant, 2012; Chen et al., 2020; Kearns and Singh, 1999; Li et al., 2023a, 2022a, 2023b; Sidford et al., 2018; Wainwright, 2019) and the offline setting (Jin et al., 2021; Li et al., 2022b; Rashidinejad et al., 2021; Shi et al., 2022; Xie et al., 2021; Yan et al., 2022; Yin et al., 2021), as well as the online setting via both regret-based and PAC-base analyses (Bai et al., 2019; Dong et al., 2019; Jafarnia-Jahromi et al., 2020; Jin et al., 2018, 2020; Li et al., 2021, 2023c; Yang et al., 2021; Zhang et al., 2020b).

**Robustness in RL.** While standard RL has achieved remarkable success, current RL algorithms still have significant drawbacks in that the learned policy could be completely off if the deployed environment is subject to perturbation, model mismatch, or other structural changes. To address these challenges, an emerging line of works begin to address robustness of RL algorithms with respect to the uncertainty or perturbation over different components of MDPs — state, action, reward, and the transition kernel; see Moos et al. (2022) for a recent review. Besides the framework of distributionally robust MDPs (RMDPs) (Iyengar, 2005) adopted by this work, to promote robustness in RL, there exist various other works including but not limited to Han et al. (2022); Qiaoben et al. (2021); Sun et al. (2021); Xiong et al. (2022); Zhang et al. (2021, 2020a) investigating the robustness w.r.t. state uncertainty, where the agent’s policy is chosen based on a perturbed observation generated from the state by adding restricted noise or adversarial attack. Besides, Tan et al. (2020); Tessler et al. (2019) considered the robustness to the uncertainty of the action, namely, the action is possibly distorted by an adversarial agent abruptly or smoothly.

**Distributionally robust RL.** Rooted in the literature of distributionally robust optimization, which has primarily been investigated in the context of supervised learning (Bertsimas et al., 2018; Blanchet and Murthy, 2019; Duchi and Namkoong, 2018; Gao, 2020; Rahimian and Mehrotra, 2019), distributionally robust dynamic programming and RMDPs have attracted considerable attention recently (Badrinath and Kalathil, 2021; Derman and Mannor, 2020; Goyal and Grand-Clement, 2022; Ho et al., 2018, 2021; Iyengar, 2005; Kaufman and Schaefer, 2013; Smirnova et al., 2019; Tamar et al., 2014; Wolff et al., 2012; Xu and Mannor, 2012). In the context of RMDPs, both empirical and theoretical studies have been widely conducted, although most prior theoretical analyses focus on planning with an exact knowledge of the uncertainty set (Iyengar, 2005; Tamar et al., 2014; Xu and Mannor, 2012), or are asymptotic in nature (Roy et al., 2017).

Resorting to the tools of high-dimensional statistics, various recent works begin to shift attention to understand the finite-sample performance of provable robust RL algorithms, under diverse data generating mechanisms and forms of the uncertainty set over the transition kernel. Besides the infinite-horizon setting, finite-sample complexity bounds for RMDPs with the TV distance and the  $\chi^2$  divergence are also developed for the finite-horizon setting in Dong et al. (2022); Xu et al. (2023). In addition, many other forms of uncertainty sets have been considered. For example, Wang and Zou (2021) considered a R-contamination uncertain set and proposed a provable robust Q-learning algorithm for the online setting with similar guarantees as standard MDPs. The KL divergence is another popular choice widely considered, where Blanchet et al. (2023); Panaganti and Kalathil (2022); Shi and Chi (2022); Wang et al. (2023); Xu et al. (2023); Yang et al. (2022); Zhou et al. (2021) investigated the sample complexity of both model-based and model-free algorithms under the simulator or offline settings. Xu et al. (2023) considered a variety of uncertainty sets including one associated with Wasserstein distance. Badrinath and Kalathil (2021) considered a general  $(s, a)$ -rectangular form of the uncertainty set and proposed a model-free algorithm for the online setting with linear function approximation to cope with large state spaces. Moreover, various other related issues

have been explored such as the iteration complexity of the policy-based methods (Kumar et al., 2023; Li et al., 2022c), and regularization-based robust RL (Yang et al., 2023).

## 6 Discussions

This work has developed improved sample complexity bounds for learning RMDPs when the uncertainty set is measured via the TV distance or the  $\chi^2$  divergence, assuming availability of a generative model. Our results have not only strengthened the prior art in both the upper and lower bounds, but have also unlocked curious insights into how the quest for distributional robustness impacts the sample complexity. As a key takeaway of this paper, RMDPs are not necessarily harder nor easier to learn than standard MDPs, as the answer depends — in a rather subtle manner — on the specific choice of the uncertainty set. For the case w.r.t. the TV distance, we have settled the minimax sample complexity for RMDPs, which is never larger than that required to learn standard MDPs. Regarding the case w.r.t. the  $\chi^2$  divergence, we have uncovered that learning RMDPs can oftentimes be provably harder than the standard MDP counterpart. All in all, our findings help raise awareness that the choice of the uncertainty set not only represents a preference in robustness, but also exerts fundamental influences upon the intrinsic statistical complexity.

Moving forward, our work opens up numerous avenues for future studies, and we point out a few below.

- *Extensions to the finite-horizon setting.* It is likely that our current analysis framework can be extended to tackle finite-horizon RMDPs, which would help complete our understanding for the tabular cases.
- *Improved analysis for the case of  $\chi^2$  divergence.* While we have settled the sample complexity of RMDPs with the TV distance, the upper and lower bounds we have developed for RMDPs w.r.t. the  $\chi^2$  divergence still differ by some polynomial factor in the effective horizon. It would be of great interest to see how to close this gap.
- *A unified theory for other families of uncertainty sets.* Our work raises an interesting question concerning how the geometry of the uncertainty sets intervenes the sample complexity. Characterizing the tight sample complexity for RMDPs under a more general family of uncertainty sets — such as using  $\ell_p$  distance or  $f$ -divergence, as well as  $s$ -rectangular sets — would be highly desirable.

## Acknowledgement

The work of L. Shi and Y. Chi is supported in part by the grants ONR N00014-19-1-2404, NSF CCF-2106778, DMS-2134080, and CNS-2148212. L. Shi is also gratefully supported by the Leo Finzi Memorial Fellowship, Wei Shen and Xuehong Zhang Presidential Fellowship, and Liang Ji-Dian Graduate Fellowship at Carnegie Mellon University. The work of Y. Wei is supported in part by the the NSF grants DMS-2147546/2015447, CAREER award DMS-2143215, CCF-2106778, and the Google Research Scholar Award. The work of Y. Chen is supported in part by the Alfred P. Sloan Research Fellowship, the Google Research Scholar Award, the AFOSR grant FA9550-22-1-0198, the ONR grant N00014-22-1-2354, and the NSF grants CCF-2221009 and CCF-1907661. The authors also acknowledge Zuxin Liu and He Wang for valuable discussions.

## References

- Agarwal, A., Kakade, S., and Yang, L. F. (2020). Model-based reinforcement learning with a generative model is minimax optimal. In *Conference on Learning Theory*, pages 67–83. PMLR.
- Azar, M., Munos, R., and Kappen, H. J. (2013a). Minimax pac bounds on the sample complexity of reinforcement learning with a generative model. *Machine learning*, 91:325–349.
- Azar, M. G., Munos, R., and Kappen, H. J. (2013b). Minimax PAC bounds on the sample complexity of reinforcement learning with a generative model. *Machine learning*, 91(3):325–349.

- Badrinath, K. P. and Kalathil, D. (2021). Robust reinforcement learning using least squares policy iteration with provable performance guarantees. In *International Conference on Machine Learning*, pages 511–520. PMLR.
- Bai, Y., Xie, T., Jiang, N., and Wang, Y.-X. (2019). Provably efficient Q-learning with low switching cost. *arXiv preprint arXiv:1905.12849*.
- Beck, C. L. and Srikant, R. (2012). Error bounds for constant step-size Q-learning. *Systems & control letters*, 61(12):1203–1208.
- Bertsimas, D., Gupta, V., and Kallus, N. (2018). Data-driven robust optimization. *Mathematical Programming*, 167(2):235–292.
- Blanchet, J., Lu, M., Zhang, T., and Zhong, H. (2023). Double pessimism is provably efficient for distributionally robust offline reinforcement learning: Generic algorithm and robust partial coverage. *arXiv preprint arXiv:2305.09659*.
- Blanchet, J. and Murthy, K. (2019). Quantifying distributional model risk via optimal transport. *Mathematics of Operations Research*, 44(2):565–600.
- Chen, Z., Maguluri, S. T., Shakkottai, S., and Shanmugam, K. (2020). Finite-sample analysis of stochastic approximation using smooth convex envelopes. *arXiv preprint arXiv:2002.00874*.
- Clavier, P., Pennec, E. L., and Geist, M. (2023). Towards minimax optimality of model-based robust reinforcement learning. *arXiv preprint arXiv:2302.05372*.
- Derman, E. and Mannor, S. (2020). Distributional robustness and regularization in reinforcement learning. *arXiv preprint arXiv:2003.02894*.
- Dong, J., Li, J., Wang, B., and Zhang, J. (2022). Online policy optimization for robust MDP. *arXiv preprint arXiv:2209.13841*.
- Dong, K., Wang, Y., Chen, X., and Wang, L. (2019). Q-learning with UCB exploration is sample efficient for infinite-horizon MDP. *arXiv preprint arXiv:1901.09311*.
- Duchi, J. and Namkoong, H. (2018). Learning models with uniform performance via distributionally robust optimization. *arXiv preprint arXiv:1810.08750*.
- Fatemi, M., Killian, T. W., Subramanian, J., and Ghassemi, M. (2021). Medical dead-ends and learning to identify high-risk states and treatments. *Advances in Neural Information Processing Systems*, 34:4856–4870.
- Gao, R. (2020). Finite-sample guarantees for wasserstein distributionally robust optimization: Breaking the curse of dimensionality. *arXiv preprint arXiv:2009.04382*.
- Goyal, V. and Grand-Clement, J. (2022). Robust markov decision processes: Beyond rectangularity. *Mathematics of Operations Research*.
- Han, S., Su, S., He, S., Han, S., Yang, H., and Miao, F. (2022). What is the solution for state adversarial multi-agent reinforcement learning? *arXiv preprint arXiv:2212.02705*.
- Ho, C. P., Petrik, M., and Wiesemann, W. (2018). Fast bellman updates for robust MDPs. In *International Conference on Machine Learning*, pages 1979–1988. PMLR.
- Ho, C. P., Petrik, M., and Wiesemann, W. (2021). Partial policy iteration for l1-robust markov decision processes. *Journal of Machine Learning Research*, 22(275):1–46.
- Iyengar, G. N. (2005). Robust dynamic programming. *Mathematics of Operations Research*, 30(2):257–280.
- Jafarnia-Jahromi, M., Wei, C.-Y., Jain, R., and Luo, H. (2020). A model-free learning algorithm for infinite-horizon average-reward MDPs with near-optimal regret. *arXiv preprint arXiv:2006.04354*.



- Jin, C., Allen-Zhu, Z., Bubeck, S., and Jordan, M. I. (2018). Is Q-learning provably efficient? In *Advances in Neural Information Processing Systems*, pages 4863–4873.
- Jin, C., Krishnamurthy, A., Simchowitz, M., and Yu, T. (2020). Reward-free exploration for reinforcement learning. In *International Conference on Machine Learning*, pages 4870–4879. PMLR.
- Jin, Y., Yang, Z., and Wang, Z. (2021). Is pessimism provably efficient for offline RL? In *International Conference on Machine Learning*, pages 5084–5096.
- Kaufman, D. L. and Schaefer, A. J. (2013). Robust modified policy iteration. *INFORMS Journal on Computing*, 25(3):396–410.
- Kearns, M. J. and Singh, S. P. (1999). Finite-sample convergence rates for Q-learning and indirect algorithms. In *Advances in neural information processing systems*, pages 996–1002.
- Klopp, O., Lounici, K., and Tsybakov, A. B. (2017). Robust matrix completion. *Probability Theory and Related Fields*, 169(1-2):523–564.
- Kober, J., Bagnell, J. A., and Peters, J. (2013). Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*, 32(11):1238–1274.
- Kumar, N., Derman, E., Geist, M., Levy, K., and Mannor, S. (2023). Policy gradient for s-rectangular robust markov decision processes. *arXiv preprint arXiv:2301.13589*.
- Lee, J., Jeon, W., Lee, B., Pineau, J., and Kim, K.-E. (2021). Optidice: Offline policy optimization via stationary distribution correction estimation. In *International Conference on Machine Learning*, pages 6120–6130. PMLR.
- Li, G., Cai, C., Chen, Y., Wei, Y., and Chi, Y. (2023a). Is Q-learning minimax optimal? a tight sample complexity analysis. *Operations Research*.
- Li, G., Chi, Y., Wei, Y., and Chen, Y. (2022a). Minimax-optimal multi-agent RL in Markov games with a generative model. *Neural Information Processing Systems*.
- Li, G., Shi, L., Chen, Y., Chi, Y., and Wei, Y. (2022b). Settling the sample complexity of model-based offline reinforcement learning. *arXiv preprint arXiv:2204.05275*.
- Li, G., Shi, L., Chen, Y., Gu, Y., and Chi, Y. (2021). Breaking the sample complexity barrier to regret-optimal model-free reinforcement learning. *Advances in Neural Information Processing Systems*, 34.
- Li, G., Wei, Y., Chi, Y., and Chen, Y. (2023b). Breaking the sample size barrier in model-based reinforcement learning with a generative model. *accepted to Operations Research*.
- Li, G., Yan, Y., Chen, Y., and Fan, J. (2023c). Minimax-optimal reward-agnostic exploration in reinforcement learning. *arXiv preprint arXiv:2304.07278*.
- Li, Y., Zhao, T., and Lan, G. (2022c). First-order policy optimization for robust markov decision process. *arXiv preprint arXiv:2209.10579*.
- Liu, S., Ngiam, K. Y., and Feng, M. (2019). Deep reinforcement learning for clinical decision support: a brief survey. *arXiv preprint arXiv:1907.09475*.
- Mahmood, A. R., Korenkevych, D., Vasan, G., Ma, W., and Bergstra, J. (2018). Benchmarking reinforcement learning algorithms on real-world robots. In *Conference on robot learning*, pages 561–591. PMLR.
- Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., and Riedmiller, M. (2013). Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*.
- Moos, J., Hansel, K., Abdulsamad, H., Stark, S., Clever, D., and Peters, J. (2022). Robust reinforcement learning: A review of foundations and recent advances. *Machine Learning and Knowledge Extraction*, 4(1):276–315.

- Nilim, A. and El Ghaoui, L. (2005). Robust control of Markov decision processes with uncertain transition matrices. *Operations Research*, 53(5):780–798.
- OpenAI (2023). Gpt-4 technical report.
- Panaganti, K. and Kalathil, D. (2022). Sample complexity of robust reinforcement learning with a generative model. In *International Conference on Artificial Intelligence and Statistics*, pages 9582–9602. PMLR.
- Qiaoben, Y., Zhou, X., Ying, C., and Zhu, J. (2021). Strategically-timed state-observation attacks on deep reinforcement learning agents. In *ICML 2021 Workshop on Adversarial Machine Learning*.
- Rahimian, H. and Mehrotra, S. (2019). Distributionally robust optimization: A review. *arXiv preprint arXiv:1908.05659*.
- Rashidinejad, P., Zhu, B., Ma, C., Jiao, J., and Russell, S. (2021). Bridging offline reinforcement learning and imitation learning: A tale of pessimism. *Neural Information Processing Systems (NeurIPS)*.
- Roy, A., Xu, H., and Pokutta, S. (2017). Reinforcement learning under model mismatch. *Advances in neural information processing systems*, 30.
- Shi, L. and Chi, Y. (2022). Distributionally robust model-based offline reinforcement learning with near-optimal sample complexity. *arXiv preprint arXiv:2208.05767*.
- Shi, L., Li, G., Wei, Y., Chen, Y., and Chi, Y. (2022). Pessimistic Q-learning for offline reinforcement learning: Towards optimal sample complexity. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162, pages 19967–20025. PMLR.
- Sidford, A., Wang, M., Wu, X., Yang, L., and Ye, Y. (2018). Near-optimal time and sample complexities for solving Markov decision processes with a generative model. In *Advances in Neural Information Processing Systems*, pages 5186–5196.
- Smirnova, E., Dohmatob, E., and Mary, J. (2019). Distributionally robust reinforcement learning. *arXiv preprint arXiv:1902.08708*.
- Sun, K., Liu, Y., Zhao, Y., Yao, H., Jui, S., and Kong, L. (2021). Exploring the training robustness of distributional reinforcement learning against noisy state observations. *arXiv preprint arXiv:2109.08776*.
- Tamar, A., Mannor, S., and Xu, H. (2014). Scaling up robust MDPs using function approximation. In *International conference on machine learning*, pages 181–189. PMLR.
- Tan, K. L., Esfandiari, Y., Lee, X. Y., and Sarkar, S. (2020). Robustifying reinforcement learning agents via action space adversarial training. In *2020 American control conference (ACC)*, pages 3959–3964. IEEE.
- Tessler, C., Efroni, Y., and Mannor, S. (2019). Action robust reinforcement learning and applications in continuous control. In *International Conference on Machine Learning*, pages 6215–6224. PMLR.
- Tsybakov, A. B. (2009). *Introduction to nonparametric estimation*, volume 11. Springer.
- Vershynin, R. (2018). *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press.
- Wainwright, M. J. (2019). Stochastic approximation with cone-contractive operators: Sharp  $\ell_\infty$ -bounds for Q-learning. *arXiv preprint arXiv:1905.06265*.
- Wang, S., Si, N., Blanchet, J., and Zhou, Z. (2023). A finite sample complexity bound for distributionally robust q-learning. *arXiv preprint arXiv:2302.13203*.
- Wang, Y. and Zou, S. (2021). Online robust reinforcement learning with model uncertainty. *Advances in Neural Information Processing Systems*, 34.
- Wiesemann, W., Kuhn, D., and Rustem, B. (2013). Robust markov decision processes. *Mathematics of Operations Research*, 38(1):153–183.

- Wolff, E. M., Topcu, U., and Murray, R. M. (2012). Robust control of uncertain markov decision processes with temporal logic specifications. In *2012 IEEE 51st IEEE Conference on Decision and Control (CDC)*, pages 3372–3379. IEEE.
- Xie, T., Jiang, N., Wang, H., Xiong, C., and Bai, Y. (2021). Policy finetuning: Bridging sample-efficient offline and online reinforcement learning. *Advances in neural information processing systems*, 34.
- Xiong, Z., Eappen, J., Zhu, H., and Jagannathan, S. (2022). Defending observation attacks in deep reinforcement learning via detection and denoising. *arXiv preprint arXiv:2206.07188*.
- Xu, H. and Mannor, S. (2012). Distributionally robust Markov decision processes. *Mathematics of Operations Research*, 37(2):288–300.
- Xu, Z., Panaganti, K., and Kalathil, D. (2023). Improved sample complexity bounds for distributionally robust reinforcement learning. *arXiv preprint arXiv:2303.02783*.
- Yan, Y., Li, G., Chen, Y., and Fan, J. (2022). The efficacy of pessimism in asynchronous Q-learning. *arXiv preprint arXiv:2203.07368*.
- Yang, K., Yang, L., and Du, S. (2021). Q-learning with logarithmic regret. In *International Conference on Artificial Intelligence and Statistics*, pages 1576–1584. PMLR.
- Yang, W., Wang, H., Kozuno, T., Jordan, S. M., and Zhang, Z. (2023). Avoiding model estimation in robust markov decision processes with a generative model. *arXiv preprint arXiv:2302.01248*.
- Yang, W., Zhang, L., and Zhang, Z. (2022). Toward theoretical understandings of robust Markov decision processes: Sample complexity and asymptotics. *The Annals of Statistics*, 50(6):3223–3248.
- Yin, M., Bai, Y., and Wang, Y.-X. (2021). Near-optimal offline reinforcement learning via double variance reduction. *arXiv preprint arXiv:2102.01748*.
- Zhang, H., Chen, H., Boning, D., and Hsieh, C.-J. (2021). Robust reinforcement learning on state observations with learned optimal adversary. *arXiv preprint arXiv:2101.08452*.
- Zhang, H., Chen, H., Xiao, C., Li, B., Liu, M., Boning, D., and Hsieh, C.-J. (2020a). Robust deep reinforcement learning against adversarial perturbations on state observations. *Advances in Neural Information Processing Systems*, 33:21024–21037.
- Zhang, Z., Zhou, Y., and Ji, X. (2020b). Almost optimal model-free reinforcement learning via reference-advantage decomposition. *Advances in Neural Information Processing Systems*, 33.
- Zhou, Z., Bai, Q., Zhou, Z., Qiu, L., Blanchet, J., and Glynn, P. (2021). Finite-sample regret bound for distributionally robust offline tabular reinforcement learning. In *International Conference on Artificial Intelligence and Statistics*, pages 3331–3339. PMLR.
- Ziegler, D. M., Stiennon, N., Wu, J., Brown, T. B., Radford, A., Amodei, D., Christiano, P., and Irving, G. (2019). Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*.

## A Preliminaries

For convenience, we introduce the notation  $[T] := \{1, \dots, T\}$  for any positive integer  $T > 0$ . Moreover, for any two vectors  $x = [x_i]_{1 \leq i \leq n}$  and  $y = [y_i]_{1 \leq i \leq n}$ , the notation  $x \leq y$  (resp.  $x \geq y$ ) means  $x_i \leq y_i$  (resp.  $x_i \geq y_i$ ) for all  $1 \leq i \leq n$ . And for any vector  $x$ , we overload the notation by letting  $x^{\circ 2} = [x(s, a)^2]_{(s,a) \in \mathcal{S} \times \mathcal{A}}$  (resp.  $x^{\circ 2} = [x(s)^2]_{s \in \mathcal{S}}$ ). With slight abuse of notation, we denote 0 (resp. 1) as the all-zero (resp. all-one) vector, and drop the subscript  $\rho$  to write  $\mathcal{U}^\sigma(\cdot) = \mathcal{U}_\rho^\sigma(\cdot)$  whenever the argument holds for all divergence  $\rho$ .

**Matrix notation.** To continue, we recall or introduce some additional matrix notation that is useful throughout the analysis.

- $P^0 \in \mathbb{R}^{SA \times S}$ : the matrix of the nominal transition kernel with  $P_{s,a}^0$  as the  $(s, a)$ -th row.
- $\hat{P}^0 \in \mathbb{R}^{SA \times S}$ : the matrix of the estimated nominal transition kernel with  $\hat{P}_{s,a}^0$  as the  $(s, a)$ -th row.
- $r \in \mathbb{R}^{SA}$ : a vector representing the reward function  $r$  (so that  $r_{(s,a)} = r(s, a)$  for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ).
- $\Pi^\pi \in \{0, 1\}^{S \times SA}$ : a projection matrix associated with a given deterministic policy  $\pi$  taking the following form

$$\Pi^\pi = \begin{pmatrix} e_{\pi(1)}^\top & 0^\top & \cdots & 0^\top \\ 0^\top & e_{\pi(2)}^\top & \cdots & 0^\top \\ \vdots & \vdots & \ddots & \vdots \\ 0^\top & 0^\top & \cdots & e_{\pi(S)}^\top \end{pmatrix}, \quad (30)$$

where  $e_{\pi(1)}^\top, e_{\pi(2)}^\top, \dots, e_{\pi(S)}^\top \in \mathbb{R}^A$  are standard basis vectors.

- $r_\pi \in \mathbb{R}^S$ : a reward vector restricted to the actions chosen by the policy  $\pi$ , namely,  $r_\pi(s) = r(s, \pi(s))$  for all  $s \in \mathcal{S}$  (or simply,  $r_\pi = \Pi^\pi r$ ).
- $\text{Var}_P(V) \in \mathbb{R}^{SA}$ : for any transition kernel  $P \in \mathbb{R}^{SA \times S}$  and vector  $V \in \mathbb{R}^S$ , we denote the  $(s, a)$ -th row of  $\text{Var}_P(V)$  as

$$\text{Var}_P(s, a) := \text{Var}_{P_{s,a}}(V). \quad (31)$$

- $P^V \in \mathbb{R}^{SA \times S}$ ,  $\hat{P}^V \in \mathbb{R}^{SA \times S}$ : the matrices representing the probability transition kernel in the uncertainty set that leads to the worst-case value for any vector  $V \in \mathbb{R}^S$ . We denote  $P_{s,a}^V$  (resp.  $\hat{P}_{s,a}^V$ ) as the  $(s, a)$ -th row of the transition matrix  $P^V$  (resp.  $\hat{P}^V$ ). In truth, the  $(s, a)$ -th rows of these transition matrices are defined as

$$P_{s,a}^V = \arg\min_{P \in \mathcal{U}^\sigma(P_{s,a}^0)} PV, \quad \text{and} \quad \hat{P}_{s,a}^V = \arg\min_{P \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} PV. \quad (32a)$$

Furthermore, we make use of the following short-hand notation:

$$P_{s,a}^{\pi,V} := P_{s,a}^{V^{\pi,\sigma}} = \arg\min_{P \in \mathcal{U}^\sigma(P_{s,a}^0)} PV^{\pi,\sigma}, \quad P_{s,a}^{\pi,\hat{V}} := P_{s,a}^{\hat{V}^{\pi,\sigma}} = \arg\min_{P \in \mathcal{U}^\sigma(P_{s,a}^0)} P\hat{V}^{\pi,\sigma}, \quad (32b)$$

$$\hat{P}_{s,a}^{\pi,V} := \hat{P}_{s,a}^{V^{\pi,\sigma}} = \arg\min_{P \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} PV^{\pi,\sigma}, \quad \hat{P}_{s,a}^{\pi,\hat{V}} := \hat{P}_{s,a}^{\hat{V}^{\pi,\sigma}} = \arg\min_{P \in \mathcal{U}^\sigma(\hat{P}_{s,a}^0)} P\hat{V}^{\pi,\sigma}. \quad (32c)$$

The corresponding probability transition matrices are denoted by  $P^{\pi,V} \in \mathbb{R}^{SA \times S}$ ,  $P^{\pi,\hat{V}} \in \mathbb{R}^{SA \times S}$ ,  $\hat{P}^{\pi,V} \in \mathbb{R}^{SA \times S}$  and  $\hat{P}^{\pi,\hat{V}} \in \mathbb{R}^{SA \times S}$ , respectively.

- $P^\pi \in \mathbb{R}^{S \times S}$ ,  $\hat{P}^\pi \in \mathbb{R}^{S \times S}$ ,  $\underline{P}^{\pi,V} \in \mathbb{R}^{S \times S}$ ,  $\underline{P}^{\pi,\hat{V}} \in \mathbb{R}^{S \times S}$ ,  $\hat{\underline{P}}^{\pi,V} \in \mathbb{R}^{S \times S}$  and  $\hat{\underline{P}}^{\pi,\hat{V}} \in \mathbb{R}^{S \times S}$ : six square probability transition matrices w.r.t. policy  $\pi$  over the states, namely

$$\begin{aligned} P^\pi &:= \Pi^\pi P^0, & \hat{P}^\pi &:= \Pi^\pi \hat{P}^0, & \underline{P}^{\pi,V} &:= \Pi^\pi P^{\pi,V}, & \underline{P}^{\pi,\hat{V}} &:= \Pi^\pi P^{\pi,\hat{V}}, \\ \hat{\underline{P}}^{\pi,V} &:= \Pi^\pi \hat{P}^{\pi,V}, & \text{and} & & \hat{\underline{P}}^{\pi,\hat{V}} &:= \Pi^\pi \hat{P}^{\pi,\hat{V}}. \end{aligned} \quad (33)$$

We denote  $P_s^\pi$  as the  $s$ -th row of the transition matrix  $P^\pi$ ; similar quantities can be defined for the other matrices as well.

## A.1 Basic facts

**Kullback-Leibler (KL) divergence.** First, for any two distributions  $P$  and  $Q$ , we denote by  $\text{KL}(P \parallel Q)$  the Kullback-Leibler (KL) divergence of  $P$  and  $Q$ . Letting  $\text{Ber}(p)$  be the Bernoulli distribution with mean  $p$ , we also introduce

$$\text{KL}(p \parallel q) := p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q} \quad \text{and} \quad \chi^2(p \parallel q) := \frac{(p-q)^2}{q} + \frac{(p-q)^2}{1-q} = \frac{(p-q)^2}{q(1-q)}, \quad (34)$$

which represent respectively the KL divergence and the  $\chi^2$  divergence of  $\text{Ber}(p)$  from  $\text{Ber}(q)$  (Tsybakov, 2009). We make note of the following useful property about the KL divergence in Tsybakov (2009, Lemma 2.7).

**Lemma 1.** *For any  $p, q \in (0, 1)$ , it holds that*

$$\text{KL}(p \parallel q) \leq \frac{(p-q)^2}{q(1-q)}. \quad (35)$$

**Variance.** For any probability vector  $P \in \mathbb{R}^{1 \times S}$  and vector  $V \in \mathbb{R}^S$ , we denote the variance

$$\text{Var}_P(V) := P(V \circ V) - (PV) \circ (PV). \quad (36)$$

The following lemma bounds the Lipschitz constant of the variance function.

**Lemma 2.** *Consider any  $0 \leq V_1, V_2 \leq \frac{1}{1-\gamma}$  obeying  $\|V_1 - V_2\|_\infty \leq x$  and any probability vector  $P \in \Delta(S)$ , one has*

$$|\text{Var}_P(V_1) - \text{Var}_P(V_2)| \leq \frac{2x}{(1-\gamma)}. \quad (37)$$

*Proof.* It is immediate to check that

$$\begin{aligned} |\text{Var}_P(V_1) - \text{Var}_P(V_2)| &= |P(V_1 \circ V_1) - (PV_1) \circ (PV_1) - P(V_2 \circ V_2) + (PV_2) \circ (PV_2)| \\ &\leq |P(V_1 \circ V_1 - V_2 \circ V_2)| + |(PV_1 + PV_2)P(V_1 - V_2)| \\ &\leq 2\|V_1 + V_2\|_\infty \|V_1 - V_2\|_\infty \leq \frac{2x}{(1-\gamma)}. \end{aligned} \quad (38)$$

where the penultimate inequality holds by the triangle inequality.  $\square$

## A.2 Properties of the robust Bellman operator

**$\gamma$ -contraction of the robust Bellman operator.** It is worth noting that the robust Bellman operator (cf. (8)) shares the nice  $\gamma$ -contraction property of the standard Bellman operator, stated as below.

**Lemma 3** ( $\gamma$ -Contraction). (Iyengar, 2005, Theorem 3.2) *For any  $\gamma \in [0, 1)$ , the robust Bellman operator  $\mathcal{T}^\sigma(\cdot)$  (cf. (8)) is a  $\gamma$ -contraction w.r.t.  $\|\cdot\|_\infty$ . Namely, for any  $Q_1, Q_2 \in \mathbb{R}^{SA}$  s.t.  $Q_1(s, a), Q_2(s, a) \in [0, \frac{1}{1-\gamma}]$  for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , one has*

$$\|\mathcal{T}^\sigma(Q_1) - \mathcal{T}^\sigma(Q_2)\|_\infty \leq \gamma \|Q_1 - Q_2\|_\infty. \quad (39)$$

Additionally,  $Q^{*,\sigma}$  is the unique fixed point of  $\mathcal{T}^\sigma(\cdot)$  obeying  $0 \leq Q^{*,\sigma}(s, a) \leq \frac{1}{1-\gamma}$  for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ .

**Dual equivalence of the robust Bellman operator.** Fortunately, the robust Bellman operator can be evaluated efficiently by resorting to its dual formulation (Iyengar, 2005). In what follows, we shall illustrate this for the two choices of the divergence  $\rho$  of interest. Before continuing, for any  $V \in \mathbb{R}^S$ , we denote  $[V]_\alpha$  as its clipped version by some non-negative value  $\alpha$ , namely,

$$[V]_\alpha(s) := \begin{cases} \alpha, & \text{if } V(s) > \alpha, \\ V(s), & \text{otherwise.} \end{cases} \quad (40)$$

- TV distance, where the uncertainty set is  $\mathcal{U}_\rho^\sigma(\hat{P}_{s,a}^0) := \mathcal{U}_{\text{TV}}^\sigma(\hat{P}_{s,a}^0) := \mathcal{U}_{\rho_{\text{TV}}}^\sigma(\hat{P}_{s,a}^0)$  w.r.t. the TV distance  $\rho = \rho_{\text{TV}}$  defined in (9). In particular, we have the following lemma due to strong duality, which is a direct consequence of [Iyengar \(2005, Lemma 4.3\)](#).

**Lemma 4** (Strong duality for TV). *Consider any probability vector  $P \in \Delta(\mathcal{S})$ , any fixed uncertainty level  $\sigma$  and the uncertainty set  $\mathcal{U}^\sigma(P) := \mathcal{U}_{\text{TV}}^\sigma(P)$ . For any vector  $V \in \mathbb{R}^S$  obeying  $V \geq 0$ , recalling the definition of  $[V]_\alpha$  in (40), one has*

$$\inf_{P \in \mathcal{U}^\sigma(P)} \mathcal{P}V = \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left\{ P[V]_\alpha - \sigma \left( \alpha - \min_{s'} [V]_\alpha(s') \right) \right\}. \quad (41)$$

In view of the above lemma, the following dual update rule is equivalent to (16) in DRVI:

$$\hat{Q}_t(s, a) = r(s, a) + \gamma \max_{\alpha \in [\min_s \hat{V}_{t-1}(s), \max_s \hat{V}_{t-1}(s)]} \left\{ \hat{P}_{s,a}^0 [\hat{V}_{t-1}]_\alpha - \sigma \left( \alpha - \min_{s'} [\hat{V}_{t-1}]_\alpha(s') \right) \right\}. \quad (42)$$

- $\chi^2$  divergence, where the uncertainty set is  $\mathcal{U}_\rho^\sigma(\hat{P}_{s,a}^0) := \mathcal{U}_{\chi^2}^\sigma(\hat{P}_{s,a}^0) := \mathcal{U}_{\rho_{\chi^2}}^\sigma(\hat{P}_{s,a}^0)$  w.r.t. the  $\chi^2$  divergence  $\rho = \rho_{\chi^2}$  defined in (10). We introduce the following lemma which directly follows from [\(Iyengar, 2005, Lemma 4.2\)](#).

**Lemma 5** (Strong duality for  $\chi^2$ ). *Consider any probability vector  $P \in \Delta(\mathcal{S})$ , any fixed uncertainty level  $\sigma$  and the uncertainty set  $\mathcal{U}^\sigma(P) := \mathcal{U}_{\chi^2}^\sigma(P)$ . For any vector  $V \in \mathbb{R}^S$  obeying  $V \geq 0$ , one has*

$$\inf_{P \in \mathcal{U}^\sigma(P)} \mathcal{P}V = \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left\{ P[V]_\alpha - \sqrt{\sigma \text{Var}_P([V]_\alpha)} \right\}, \quad (43)$$

where  $\text{Var}_P(\cdot)$  is defined as (36).

In view of the above lemma, the update rule (16) in DRVI can be equivalently written as:

$$\hat{Q}_t(s, a) = r(s, a) + \gamma \max_{\alpha \in [\min_s \hat{V}_{t-1}(s), \max_s \hat{V}_{t-1}(s)]} \left\{ \hat{P}_{s,a}^0 [\hat{V}_{t-1}]_\alpha - \sqrt{\sigma \text{Var}_{\hat{P}_{s,a}^0}([\hat{V}_{t-1}]_\alpha)} \right\}. \quad (44)$$

The proofs of Lemma 4 and Lemma 5 are provided as follows.

*Proof of Lemma 4.* To begin with, applying [\(Iyengar, 2005, Lemma 4.3\)](#), the term of interest obeys

$$\inf_{P \in \mathcal{U}^\sigma(P)} \mathcal{P}V = \max_{\mu \in \mathbb{R}^S, \mu \geq 0} \left\{ P(V - \mu) - \sigma \left( \max_{s'} \{V(s') - \mu(s')\} - \min_{s'} \{V(s') - \mu(s')\} \right) \right\}, \quad (45)$$

where  $\mu(s')$  represents the  $s'$ -th entry of  $\mu \in \mathbb{R}^S$ . Denoting  $\mu^*$  as the optimal dual solution, taking  $\alpha = \max_{s'} \{V(s') - \mu^*(s')\}$ , it is easily verified that  $\mu^*$  obeys

$$\mu^*(s) = \begin{cases} V(s) - \alpha, & \text{if } V(s) > \alpha \\ 0, & \text{otherwise.} \end{cases} \quad (46)$$

Therefore, (45) can be solved by optimizing  $\alpha$  as below [\(Iyengar, 2005, Lemma 4.3\)](#):

$$\inf_{P \in \mathcal{U}^\sigma(P)} \mathcal{P}V = \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left\{ P[V]_\alpha - \sigma \left( \alpha - \min_{s'} [V]_\alpha(s') \right) \right\}. \quad (47)$$

□

*Proof of Lemma 5.* Due to strong duality [\(Iyengar, 2005, Lemma 4.2\)](#), it holds that

$$\inf_{P \in \mathcal{U}^\sigma(P)} \mathcal{P}V = \max_{\mu \in \mathbb{R}^S, \mu \geq 0} \left\{ P(V - \mu) - \sqrt{\sigma \text{Var}_P(V - \mu)} \right\}, \quad (48)$$

and the optimal  $\mu^*$  obeys

$$\mu^*(s) = \begin{cases} V(s) - \alpha, & \text{if } V(s) > \alpha \\ 0, & \text{otherwise.} \end{cases} \quad (49)$$

for some  $\alpha \in [\min_s V(s), \max_s V(s)]$ . As a result, solving (48) is equivalent to optimizing the scalar  $\alpha$  as below:

$$\inf_{\mathcal{P} \in \mathcal{U}^\sigma(P)} \mathcal{P}V = \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left\{ P[V]_\alpha - \sqrt{\sigma \text{Var}_P([V]_\alpha)} \right\}. \quad (50)$$

□

### A.3 Additional facts of the empirical robust MDP

**Bellman equations of the empirical robust MDP  $\widehat{\mathcal{M}}_{\text{rob}}$ .** To begin with, recall that the empirical robust MDP  $\widehat{\mathcal{M}}_{\text{rob}} = \{\mathcal{S}, \mathcal{A}, \gamma, \mathcal{U}^\sigma(\widehat{P}^0), r\}$  based on the estimated nominal distribution  $\widehat{P}^0$  constructed in (13) and its corresponding robust value function (resp. robust Q-function)  $\widehat{V}^{\pi, \sigma}$  (resp.  $\widehat{Q}^{\pi, \sigma}$ ).

Note that  $\widehat{Q}^{\pi, \sigma}$  is the unique fixed point of  $\widehat{\mathcal{T}}^\sigma(\cdot)$  (see Lemma 3), the empirical robust Bellman operator constructed using  $\widehat{P}^0$ . Moreover, similar to (7), for  $\widehat{\mathcal{M}}_{\text{rob}}$ , the Bellman's optimality principle gives the following *robust Bellman consistency equation* (resp. *robust Bellman optimality equation*):

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad \widehat{Q}^{\pi, \sigma}(s, a) = r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\widehat{P}_{s,a}^0)} \mathcal{P}\widehat{V}^{\pi, \sigma}, \quad (51a)$$

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad \widehat{Q}^{\star, \sigma}(s, a) = r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\widehat{P}_{s,a}^0)} \mathcal{P}\widehat{V}^{\star, \sigma}. \quad (51b)$$

With these in mind, combined with the matrix notation, for any policy  $\pi$ , we can write the robust Bellman consistency equations as

$$Q^{\pi, \sigma} = r + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P^0)} \mathcal{P}V^{\pi, \sigma} \quad \text{and} \quad \widehat{Q}^{\pi, \sigma} = r + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\widehat{P}^0)} \mathcal{P}\widehat{V}^{\pi, \sigma}, \quad (52)$$

which leads to

$$\begin{aligned} V^{\pi, \sigma} &= r_\pi + \gamma \Pi^\pi \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P^0)} \mathcal{P}V^{\pi, \sigma} \stackrel{(i)}{=} r_\pi + \gamma \underline{P}^{\pi, V} V^{\pi, \sigma}, \\ \widehat{V}^{\pi, \sigma} &= r_\pi + \gamma \Pi^\pi \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\widehat{P}^0)} \mathcal{P}\widehat{V}^{\pi, \sigma} \stackrel{(ii)}{=} r_\pi + \gamma \widehat{\underline{P}}^{\pi, \widehat{V}} \widehat{V}^{\pi, \sigma}, \end{aligned} \quad (53)$$

where (i) and (ii) holds by the definitions in (30), (32) and (33).

Encouragingly, the above property of the robust Bellman operator ensures the fast convergence of DRVI. We collect this consequence in the following lemma.

**Lemma 6.** *Let  $\widehat{Q}_0 = 0$ . The iterates  $\{\widehat{Q}_t\}, \{\widehat{V}_t\}$  of DRVI obey*

$$\forall t \geq 0: \quad \|\widehat{Q}_t - \widehat{Q}^{\star, \sigma}\|_\infty \leq \frac{\gamma^t}{1 - \gamma} \quad \text{and} \quad \|\widehat{V}_t - \widehat{V}^{\star, \sigma}\|_\infty \leq \frac{\gamma^t}{1 - \gamma}. \quad (54)$$

Furthermore, the output policy  $\widehat{\pi}$  obeys

$$\|\widehat{V}^{\star, \sigma} - \widehat{V}^{\widehat{\pi}, \sigma}\|_\infty \leq \frac{2\gamma\varepsilon_{\text{opt}}}{1 - \gamma}, \quad \text{where} \quad \|\widehat{V}^{\star, \sigma} - \widehat{V}_{T-1}\|_\infty =: \varepsilon_{\text{opt}}. \quad (55)$$

*Proof of Lemma 6.* Applying the  $\gamma$ -contraction property in Lemma 3 directly yields that for any  $t \geq 0$ ,

$$\|\widehat{Q}_t - \widehat{Q}^{\star, \sigma}\|_\infty = \|\widehat{\mathcal{T}}^\sigma(\widehat{Q}_{t-1}) - \widehat{\mathcal{T}}^\sigma(\widehat{Q}^{\star, \sigma})\|_\infty \leq \gamma \|\widehat{Q}_{t-1} - \widehat{Q}^{\star, \sigma}\|_\infty$$



$$\leq \dots \leq \gamma^t \|\hat{Q}_0 - \hat{Q}^{\star,\sigma}\|_\infty = \gamma^t \|\hat{Q}^{\star,\sigma}\|_\infty \leq \frac{\gamma^t}{1-\gamma},$$

where the last inequality holds by the fact  $\|\hat{Q}^{\star,\sigma}\|_\infty \leq \frac{1}{1-\gamma}$  (see Lemma 3). In addition,

$$\|\hat{V}_t - \hat{V}^{\star,\sigma}\|_\infty = \max_{s \in \mathcal{S}} \left\| \max_{a \in \mathcal{A}} \hat{Q}_t(s, a) - \max_{a \in \mathcal{A}} \hat{Q}^{\star,\sigma}(s, a) \right\|_\infty \leq \|\hat{Q}_t - \hat{Q}^{\star,\sigma}\|_\infty \leq \frac{\gamma^t}{1-\gamma},$$

where the penultimate inequality holds by the maximum operator is 1-Lipschitz. This completes the proof of (54).

We now move to establish (55). Note that there exists at least one state  $s_0 \in \mathcal{S}$  that is associated with the maximum of the value gap, i.e.,

$$\|\hat{V}^{\star,\sigma} - \hat{V}^{\hat{\pi},\sigma}\|_\infty = \hat{V}^{\star,\sigma}(s_0) - \hat{V}^{\hat{\pi},\sigma}(s_0) \geq \hat{V}^{\star,\sigma}(s) - \hat{V}^{\hat{\pi},\sigma}(s), \quad \forall s \in \mathcal{S}.$$

Recall  $\hat{\pi}^*$  is the optimal robust policy for the empirical RMDP  $\widehat{\mathcal{M}}_{\text{rob}}$ . For convenience, we denote  $a_1 = \hat{\pi}^*(s_0)$  and  $a_2 = \hat{\pi}(s_0)$ . Then, since  $\hat{\pi}$  is the greedy policy w.r.t.  $\hat{Q}_T$ , one has

$$r(s_0, a_1) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s_0, a_1}^0)} \mathcal{P} \hat{V}_{T-1} = \hat{Q}_T(s_0, a_1) \leq \hat{Q}_T(s_0, a_2) = r(s_0, a_2) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s_0, a_2}^0)} \mathcal{P} \hat{V}_{T-1}. \quad (56)$$

Recalling the notation in (32), the above fact and (55) altogether yield

$$\begin{aligned} r(s_0, a_1) + \gamma \hat{P}_{s_0, a_1}^{\hat{V}_{T-1}} (\hat{V}^{\star,\sigma} - \varepsilon_{\text{opt}} 1) &\leq r(s_0, a_1) + \gamma \hat{P}_{s_0, a_1}^{\hat{V}_{T-1}} \hat{V}_{T-1} \\ &\leq r(s_0, a_2) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s_0, a_2}^0)} \mathcal{P} \hat{V}_{T-1} \\ &\stackrel{(i)}{\leq} r(s_0, a_2) + \gamma \hat{P}_{s_0, a_2}^{\hat{V}^{\hat{\pi},\sigma}} \hat{V}_{T-1} \\ &\leq r(s_0, a_2) + \gamma \hat{P}_{s_0, a_2}^{\hat{V}^{\hat{\pi},\sigma}} (\hat{V}^{\star,\sigma} + \varepsilon_{\text{opt}} 1), \end{aligned} \quad (57)$$

where (i) follows from the optimality criteria. The term of interest can be controlled as

$$\begin{aligned} \|\hat{V}^{\star,\sigma} - \hat{V}^{\hat{\pi},\sigma}\|_\infty &= \hat{V}^{\star,\sigma}(s_0) - \hat{V}^{\hat{\pi},\sigma}(s_0) \\ &= r(s_0, a_1) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s_0, a_1}^0)} \mathcal{P} \hat{V}^{\star,\sigma} - \left( r(s_0, a_2) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s_0, a_2}^0)} \mathcal{P} \hat{V}^{\hat{\pi},\sigma} \right) \\ &= r(s_0, a_1) - r(s_0, a_2) + \gamma \left( \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s_0, a_1}^0)} \mathcal{P} \hat{V}^{\star,\sigma} - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s_0, a_2}^0)} \mathcal{P} \hat{V}^{\hat{\pi},\sigma} \right) \\ &\stackrel{(i)}{\leq} 2\gamma \varepsilon_{\text{opt}} + \gamma \left( \hat{P}_{s_0, a_2}^{\hat{V}^{\hat{\pi},\sigma}} \hat{V}^{\star,\sigma} - \hat{P}_{s_0, a_1}^{\hat{V}_{T-1}} \hat{V}^{\star,\sigma} + \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s_0, a_1}^0)} \mathcal{P} \hat{V}^{\star,\sigma} - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s_0, a_2}^0)} \mathcal{P} \hat{V}^{\hat{\pi},\sigma} \right) \\ &= 2\gamma \varepsilon_{\text{opt}} + \gamma \left( \hat{P}_{s_0, a_2}^{\hat{V}^{\hat{\pi},\sigma}} \hat{V}^{\star,\sigma} - \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s_0, a_2}^0)} \mathcal{P} \hat{V}^{\hat{\pi},\sigma} \right) + \gamma \left( \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s_0, a_1}^0)} \mathcal{P} \hat{V}^{\star,\sigma} - \hat{P}_{s_0, a_1}^{\hat{V}_{T-1}} \hat{V}^{\star,\sigma} \right) \\ &\stackrel{(ii)}{\leq} 2\gamma \varepsilon_{\text{opt}} + \gamma \hat{P}_{s_0, a_2}^{\hat{V}^{\hat{\pi},\sigma}} (\hat{V}^{\star,\sigma} - \hat{V}^{\hat{\pi},\sigma}) + \gamma \left( \hat{P}_{s_0, a_1}^{\hat{V}_{T-1}} \hat{V}^{\star,\sigma} - \hat{P}_{s_0, a_1}^{\hat{V}_{T-1}} \hat{V}^{\star,\sigma} \right) \\ &\leq 2\gamma \varepsilon_{\text{opt}} + \gamma \|\hat{V}^{\star,\sigma} - \hat{V}^{\hat{\pi},\sigma}\|_\infty, \end{aligned} \quad (58)$$

where (i) holds by plugging in (57), and (ii) follows from  $\inf_{\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s_0, a_1}^0)} \mathcal{P} \hat{V}^{\star,\sigma} \leq \mathcal{P} \hat{V}^{\star,\sigma}$  for any  $\mathcal{P} \in \mathcal{U}^\sigma(\hat{P}_{s_0, a_1}^0)$ . Rearranging (58) leads to

$$\|\hat{V}^{\star,\sigma} - \hat{V}^{\hat{\pi},\sigma}\|_\infty \leq \frac{2\gamma \varepsilon_{\text{opt}}}{1-\gamma}.$$

□

## B Proof of the upper bound with TV distance: Theorem 1

Throughout this section, for any transition kernel  $P$ , the uncertainty set is taken as (see (9))

$$\mathcal{U}^\sigma(P) := \mathcal{U}_{\text{TV}}^\sigma(P) = \otimes \mathcal{U}_{\text{TV}}^\sigma(P_{s,a}), \quad \mathcal{U}_{\text{TV}}^\sigma(P_{s,a}) := \left\{ P'_{s,a} \in \Delta(\mathcal{S}) : \frac{1}{2} \|P'_{s,a} - P_{s,a}\|_1 \leq \sigma \right\}. \quad (59)$$

### B.1 Technical lemmas

We begin with a key lemma concerning the dynamic range of the robust value function  $V^{\pi,\sigma}$  (cf. (5)), which produces tighter control when  $\sigma$  is large; the proof is deferred to Appendix B.3.1.

**Lemma 7.** *For any nominal transition kernel  $P \in \mathbb{R}^{\mathcal{S} \times \mathcal{A} \times \mathcal{S}}$ , any fixed uncertainty level  $\sigma$ , and any policy  $\pi$ , its corresponding robust value function  $V^{\pi,\sigma}$  (cf. (5)) satisfies*

$$\max_{s \in \mathcal{S}} V^{\pi,\sigma}(s) - \min_{s \in \mathcal{S}} V^{\pi,\sigma}(s) \leq \frac{1}{\gamma \max\{1 - \gamma, \sigma\}}.$$

Next, we introduce the following lemma, whose proof is postponed in Appendix B.3.2.

**Lemma 8.** *Consider an MDP with transition kernel matrix  $P$  and reward function  $0 \leq r \leq 1$ . For any policy  $\pi$  and its associated state transition matrix  $P_\pi := \Pi^\pi P$  and value function  $0 \leq V^{\pi,P} \leq \frac{1}{1-\gamma}$  (cf. (1)), one has*

$$(I - \gamma P_\pi)^{-1} \sqrt{\text{Var}_{P_\pi}(V^{\pi,P})} \leq \sqrt{\frac{8(\max_s V^{\pi,P}(s) - \min_s V^{\pi,P}(s))}{\gamma^2(1-\gamma)^2}} 1.$$

### B.2 Proof of Theorem 1

The main proof idea of Theorem 1 is similar to that of Agarwal et al. (2020) and Li et al. (2023b) while the argument needs essential adjustments in order to adapt to the robustness setting. Before proceeding, applying Lemma 6 yields that for any  $\varepsilon_{\text{opt}} > 0$ , as long as  $T \geq \log(\frac{1}{(1-\gamma)\varepsilon_{\text{opt}}})$ , one has

$$\|\widehat{V}^{\star,\sigma} - \widehat{V}^{\widehat{\pi},\sigma}\|_\infty \leq \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma}, \quad (60)$$

allowing us to justify the more general statement in Remark 1. To control the performance gap  $\|V^{\star,\sigma} - V^{\widehat{\pi},\sigma}\|_\infty$ , the proof is divided into several key steps.

**Step 1: decomposing the error.** Recall the optimal robust policy  $\pi^\star$  w.r.t.  $\mathcal{M}_{\text{rob}}$  and the optimal robust policy  $\widehat{\pi}^\star$ , the optimal robust value function  $\widehat{V}^{\star,\sigma}$  (resp. robust value function  $\widehat{Q}^{\pi,\sigma}$ ) w.r.t.  $\widehat{\mathcal{M}}_{\text{rob}}$ . The term of interest  $V^{\star,\sigma} - V^{\widehat{\pi},\sigma}$  can be decomposed as

$$\begin{aligned} V^{\star,\sigma} - V^{\widehat{\pi},\sigma} &= (V^{\pi^\star,\sigma} - \widehat{V}^{\pi^\star,\sigma}) + (\widehat{V}^{\pi^\star,\sigma} - \widehat{V}^{\widehat{\pi}^\star,\sigma}) + (\widehat{V}^{\widehat{\pi}^\star,\sigma} - \widehat{V}^{\widehat{\pi},\sigma}) + (\widehat{V}^{\widehat{\pi},\sigma} - V^{\widehat{\pi},\sigma}) \\ &\stackrel{(i)}{\leq} (V^{\pi^\star,\sigma} - \widehat{V}^{\pi^\star,\sigma}) + (\widehat{V}^{\widehat{\pi}^\star,\sigma} - \widehat{V}^{\widehat{\pi},\sigma}) + (\widehat{V}^{\widehat{\pi},\sigma} - V^{\widehat{\pi},\sigma}) \\ &\stackrel{(ii)}{\leq} (V^{\pi^\star,\sigma} - \widehat{V}^{\pi^\star,\sigma}) + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} 1 + (\widehat{V}^{\widehat{\pi},\sigma} - V^{\widehat{\pi},\sigma}) \end{aligned} \quad (61)$$

where (i) holds by  $\widehat{V}^{\pi^\star,\sigma} - \widehat{V}^{\widehat{\pi}^\star,\sigma} \leq 0$  since  $\widehat{\pi}^\star$  is the robust optimal policy for  $\widehat{\mathcal{M}}_{\text{rob}}$ , and (ii) comes from the fact in (60).

To control the two important terms in (61), we first consider a more general term  $\widehat{V}^{\pi,\sigma} - V^{\pi,\sigma}$  for any policy  $\pi$ . Towards this, plugging in (53) yields

$$\widehat{V}^{\pi,\sigma} - V^{\pi,\sigma} = r_\pi + \gamma \widehat{P}^{\pi,\widehat{V}} \widehat{V}^{\pi,\sigma} - (r_\pi + \gamma \underline{P}^{\pi,V} V^{\pi,\sigma})$$

$$\begin{aligned}
&= \left( \gamma \hat{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} - \gamma \underline{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} \right) + \left( \gamma \underline{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} - \gamma \underline{P}^{\pi, V} V^{\pi, \sigma} \right) \\
&\stackrel{(i)}{\leq} \gamma \left( \underline{P}^{\pi, V} \hat{V}^{\pi, \sigma} - \underline{P}^{\pi, V} V^{\pi, \sigma} \right) + \left( \gamma \hat{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} - \gamma \underline{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} \right),
\end{aligned}$$

where (i) holds by observing

$$\underline{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} \leq \underline{P}^{\pi, V} \hat{V}^{\pi, \sigma}$$

due to the optimality of  $\underline{P}^{\pi, \hat{V}}$  (cf. (32)). Rearranging terms leads to

$$\hat{V}^{\pi, \sigma} - V^{\pi, \sigma} \leq \gamma \left( I - \gamma \underline{P}^{\pi, V} \right)^{-1} \left( \hat{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} - \underline{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} \right). \quad (62)$$

Similarly, we can also deduce

$$\begin{aligned}
\hat{V}^{\pi, \sigma} - V^{\pi, \sigma} &= r_{\pi} + \gamma \hat{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} - (r_{\pi} + \gamma \underline{P}^{\pi, V} V^{\pi, \sigma}) \\
&= \left( \gamma \hat{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} - \gamma \underline{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} \right) + \left( \gamma \underline{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} - \gamma \underline{P}^{\pi, V} V^{\pi, \sigma} \right) \\
&\geq \gamma \left( \underline{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} - \underline{P}^{\pi, \hat{V}} V^{\pi, \sigma} \right) + \left( \gamma \hat{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} - \gamma \underline{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} \right) \\
&\geq \gamma \left( I - \gamma \underline{P}^{\pi, \hat{V}} \right)^{-1} \left( \hat{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} - \underline{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} \right).
\end{aligned} \quad (63)$$

Combining (62) and (63), we arrive at

$$\begin{aligned}
\|\hat{V}^{\pi, \sigma} - V^{\pi, \sigma}\|_{\infty} &\leq \gamma \max \left\{ \left\| \left( I - \gamma \underline{P}^{\pi, V} \right)^{-1} \left( \hat{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} - \underline{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} \right) \right\|_{\infty}, \right. \\
&\quad \left. \left\| \left( I - \gamma \underline{P}^{\pi, \hat{V}} \right)^{-1} \left( \hat{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} - \underline{P}^{\pi, \hat{V}} \hat{V}^{\pi, \sigma} \right) \right\|_{\infty} \right\}.
\end{aligned} \quad (64)$$

By decomposing the error in a symmetric way, we can similarly obtain

$$\begin{aligned}
\|\hat{V}^{\pi, \sigma} - V^{\pi, \sigma}\|_{\infty} &\leq \gamma \max \left\{ \left\| \left( I - \gamma \hat{P}^{\pi, V} \right)^{-1} \left( \hat{P}^{\pi, V} V^{\pi, \sigma} - \underline{P}^{\pi, V} V^{\pi, \sigma} \right) \right\|_{\infty}, \right. \\
&\quad \left. \left\| \left( I - \gamma \hat{P}^{\pi, \hat{V}} \right)^{-1} \left( \hat{P}^{\pi, V} V^{\pi, \sigma} - \underline{P}^{\pi, V} V^{\pi, \sigma} \right) \right\|_{\infty} \right\}.
\end{aligned} \quad (65)$$

With the above facts in mind, we are ready to control the two terms  $\|\hat{V}^{\pi^*, \sigma} - V^{\pi^*, \sigma}\|_{\infty}$  and  $\|\hat{V}^{\hat{\pi}, \sigma} - V^{\hat{\pi}, \sigma}\|_{\infty}$  in (61) separately. More specifically, taking  $\pi = \pi^*$ , applying (65) leads to

$$\begin{aligned}
\|\hat{V}^{\pi^*, \sigma} - V^{\pi^*, \sigma}\|_{\infty} &\leq \gamma \max \left\{ \left\| \left( I - \gamma \hat{P}^{\pi^*, V} \right)^{-1} \left( \hat{P}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right) \right\|_{\infty}, \right. \\
&\quad \left. \left\| \left( I - \gamma \hat{P}^{\pi^*, \hat{V}} \right)^{-1} \left( \hat{P}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right) \right\|_{\infty} \right\}.
\end{aligned} \quad (66)$$

Similarly, taking  $\pi = \hat{\pi}$ , applying (64) leads to

$$\begin{aligned}
\|\hat{V}^{\hat{\pi}, \sigma} - V^{\hat{\pi}, \sigma}\|_{\infty} &\leq \gamma \max \left\{ \left\| \left( I - \gamma \underline{P}^{\hat{\pi}, \hat{V}} \right)^{-1} \left( \hat{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} \right) \right\|_{\infty}, \right. \\
&\quad \left. \left\| \left( I - \gamma \underline{P}^{\hat{\pi}, V} \right)^{-1} \left( \hat{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} \right) \right\|_{\infty} \right\}.
\end{aligned} \quad (67)$$

**Step 2: controlling  $\|\hat{V}^{\pi^*, \sigma} - V^{\pi^*, \sigma}\|_{\infty}$ : bounding the first term in (66).** To control the two terms in (66), we first introduce the following lemma whose proof is postponed to Appendix B.3.3.

**Lemma 9.** Consider any  $\delta \in (0, 1)$ . Setting  $N \geq \log(\frac{18SAN}{\delta})$ , with probability at least  $1 - \delta$ , one has

$$\begin{aligned} \left| \hat{\underline{P}}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right| &\leq 2 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \sqrt{\text{Var}_{P^{\pi^*}}(V^{\star, \sigma})} + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} 1 \\ &\leq 3 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} 1, \end{aligned} \quad (68)$$

where  $\text{Var}_{P^{\pi^*}}(V^{\star, \sigma})$  is defined in (31).

Armed with the above lemma, now we control the first term on the right hand side of (66) as follows:

$$\begin{aligned} &\left( I - \gamma \hat{\underline{P}}^{\pi^*, V} \right)^{-1} \left( \hat{\underline{P}}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right) \\ &\stackrel{(i)}{\leq} \left( I - \gamma \hat{\underline{P}}^{\pi^*, V} \right)^{-1} \left\| \hat{\underline{P}}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right\|_{\infty} \\ &\stackrel{(ii)}{\leq} \left( I - \gamma \hat{\underline{P}}^{\pi^*, V} \right)^{-1} \left( 2 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \sqrt{\text{Var}_{P^{\pi^*}}(V^{\star, \sigma})} + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} 1 \right) \\ &\leq \underbrace{\frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} \left( I - \gamma \hat{\underline{P}}^{\pi^*, V} \right)^{-1} 1 + 2 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left( I - \gamma \hat{\underline{P}}^{\pi^*, V} \right)^{-1} \sqrt{\text{Var}_{\hat{\underline{P}}^{\pi^*, V}}(V^{\star, \sigma})}}_{=: \mathcal{C}_1} \\ &\quad + \underbrace{2 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left( I - \gamma \hat{\underline{P}}^{\pi^*, V} \right)^{-1} \sqrt{\left| \text{Var}_{\hat{\underline{P}}^{\pi^*}}(V^{\star, \sigma}) - \text{Var}_{\hat{\underline{P}}^{\pi^*, V}}(V^{\star, \sigma}) \right|}}_{=: \mathcal{C}_2} \\ &\quad + \underbrace{2 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left( I - \gamma \hat{\underline{P}}^{\pi^*, V} \right)^{-1} \left( \sqrt{\text{Var}_{P^{\pi^*}}(V^{\star, \sigma})} - \sqrt{\text{Var}_{\hat{\underline{P}}^{\pi^*}}(V^{\star, \sigma})} \right)}_{=: \mathcal{C}_3}, \end{aligned} \quad (69)$$

where (i) holds by  $\left( I - \gamma \hat{\underline{P}}^{\pi^*, V} \right)^{-1} \geq 0$ , (ii) follows from Lemma 9, and the last inequality arise from

$$\begin{aligned} \sqrt{\text{Var}_{P^{\pi^*}}(V^{\star, \sigma})} &= \left( \sqrt{\text{Var}_{P^{\pi^*}}(V^{\star, \sigma})} - \sqrt{\text{Var}_{\hat{\underline{P}}^{\pi^*}}(V^{\star, \sigma})} \right) + \sqrt{\text{Var}_{\hat{\underline{P}}^{\pi^*}}(V^{\star, \sigma})} \\ &\leq \left( \sqrt{\text{Var}_{P^{\pi^*}}(V^{\star, \sigma})} - \sqrt{\text{Var}_{\hat{\underline{P}}^{\pi^*}}(V^{\star, \sigma})} \right) + \sqrt{\left| \text{Var}_{\hat{\underline{P}}^{\pi^*}}(V^{\star, \sigma}) - \text{Var}_{\hat{\underline{P}}^{\pi^*, V}}(V^{\star, \sigma}) \right|} + \sqrt{\text{Var}_{\hat{\underline{P}}^{\pi^*, V}}(V^{\star, \sigma})} \end{aligned}$$

by applying the triangle inequality.

To continue, observing that each row of  $\hat{\underline{P}}^{\pi^*, V}$  is a probability distribution obeying that the sum is 1, we arrive at

$$\left( I - \gamma \hat{\underline{P}}^{\pi^*, V} \right)^{-1} 1 = \left( I + \sum_{t=1}^{\infty} \gamma^t \left( \hat{\underline{P}}^{\pi^*, V} \right)^t \right) 1 = \frac{1}{1-\gamma} 1. \quad (70)$$

Armed with this fact, we shall control the other three terms  $\mathcal{C}_1, \mathcal{C}_2, \mathcal{C}_3$  in (69) separately.

- Consider  $\mathcal{C}_1$ . We first introduce the following lemma, whose proof is postponed to Appendix B.3.4.

**Lemma 10.** Consider any  $\delta \in (0, 1)$ . With probability at least  $1 - \delta$ , one has

$$\left( I - \gamma \hat{\underline{P}}^{\pi^*, V} \right)^{-1} \sqrt{\text{Var}_{\hat{\underline{P}}^{\pi^*, V}}(V^{\star, \sigma})} \leq 4 \sqrt{\frac{\left( 1 + \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \right)}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, \sigma\}}} 1 \leq 4 \sqrt{\frac{\left( 1 + \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \right)}{\gamma^3 (1-\gamma)^3}} 1.$$

Applying Lemma 10 and inserting back to (69) leads to

$$\begin{aligned} \mathcal{C}_1 &= 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \hat{P}^{\pi^*, V}\right)^{-1} \sqrt{\text{Var}_{\hat{P}^{\pi^*, V}}(V^{*, \sigma})} \\ &\leq 8\sqrt{\frac{\log(\frac{18SAN}{\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} \left(1 + \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}}\right) 1. \end{aligned} \quad (71)$$

- Consider  $\mathcal{C}_2$ . First, denote  $V' := V^{*, \sigma} - \min_{s' \in \mathcal{S}} V^{*, \sigma}(s')1$ , by Lemma 7, it follows that

$$0 \leq V' \leq \frac{1}{\gamma \max\{1-\gamma, \sigma\}} 1. \quad (72)$$

Then, we have for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , and  $P_{s,a} \in \Delta(\mathcal{S})$ , and  $\tilde{P}_{s,a} \in \mathcal{U}^\sigma(P_{s,a})$ :

$$\begin{aligned} |\text{Var}_{\tilde{P}_{s,a}}(V^{*, \sigma}) - \text{Var}_{P_{s,a}}(V^{*, \sigma})| &= |\text{Var}_{\tilde{P}_{s,a}}(V') - \text{Var}_{P_{s,a}}(V')| \\ &\leq \|\tilde{P}_{s,a} - P_{s,a}\|_1 \|V'\|_\infty^2 \\ &\leq \frac{2\sigma}{\gamma^2(\max\{1-\gamma, \sigma\})^2} 1 \leq \frac{2}{\gamma^2 \max\{1-\gamma, \sigma\}} 1. \end{aligned} \quad (73)$$

Applying the above relation we obtain

$$\begin{aligned} \mathcal{C}_2 &= 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \hat{P}^{\pi^*, V}\right)^{-1} \sqrt{|\text{Var}_{\hat{P}^{\pi^*}}(V^{*, \sigma}) - \text{Var}_{\hat{P}^{\pi^*, V}}(V^{*, \sigma})|} \\ &= 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \hat{P}^{\pi^*, V}\right)^{-1} \sqrt{|\Pi^{\pi^*}(\text{Var}_{\hat{P}_0}(V^{*, \sigma}) - \text{Var}_{\hat{P}^{\pi^*, V}}(V^{*, \sigma}))|} \\ &\leq 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \hat{P}^{\pi^*, V}\right)^{-1} \sqrt{\|\text{Var}_{\hat{P}_0}(V^{*, \sigma}) - \text{Var}_{\hat{P}^{\pi^*, V}}(V^{*, \sigma})\|_\infty} 1 \\ &\leq 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \hat{P}^{\pi^*, V}\right)^{-1} \sqrt{\frac{2}{\gamma^2 \max\{1-\gamma, \sigma\}}} 1 = 2\sqrt{\frac{2 \log(\frac{18SAN}{\delta})}{\gamma^2(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1, \end{aligned} \quad (74)$$

where the last equality uses  $\left(I - \gamma \hat{P}^{\pi^*, V}\right)^{-1} 1 = \frac{1}{1-\gamma}$  (cf. (70)).

- Consider  $\mathcal{C}_3$ . The following lemma plays an important role.

**Lemma 11.** (*Panaganti and Kalathil, 2022, Lemma 6*) Consider any  $\delta \in (0, 1)$ . For any fixed policy  $\pi$  and fixed value vector  $V \in \mathbb{R}^S$ , one has with probability at least  $1 - \delta$ ,

$$\left| \sqrt{\text{Var}_{\hat{P}^\pi}(V)} - \sqrt{\text{Var}_{P^\pi}(V)} \right| \leq \sqrt{\frac{2\|V\|_\infty^2 \log(\frac{2SA}{\delta})}{N}} 1.$$

Applying Lemma 11 with  $\pi = \pi^*$  and  $V = V^{*, \sigma}$  leads to

$$\sqrt{\text{Var}_{P^{\pi^*}}(V^{*, \sigma})} - \sqrt{\text{Var}_{\hat{P}^{\pi^*}}(V^{*, \sigma})} \leq \sqrt{\frac{2\|V^{*, \sigma}\|_\infty^2 \log(\frac{2SA}{\delta})}{N}} 1,$$

which can be plugged in (69) to verify

$$\mathcal{C}_3 = 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \left(I - \gamma \hat{P}^{\pi^*, V}\right)^{-1} \left( \sqrt{\text{Var}_{P^{\pi^*}}(V^{*, \sigma})} - \sqrt{\text{Var}_{\hat{P}^{\pi^*}}(V^{*, \sigma})} \right)$$

$$\leq \frac{4}{(1-\gamma)} \frac{\log(\frac{SAN}{\delta}) \|V^{\star, \sigma}\|_{\infty}}{N} 1 \leq \frac{4 \log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N} 1, \quad (75)$$

where the last line uses  $(I - \gamma \hat{P}^{\pi^*, V})^{-1} 1 = \frac{1}{1-\gamma}$  (cf. (70)).

Finally, inserting the results of  $\mathcal{C}_1$  in (71),  $\mathcal{C}_2$  in (74),  $\mathcal{C}_3$  in (75), and (70) back into (69) gives

$$\begin{aligned} (I - \gamma \hat{P}^{\pi^*, V})^{-1} (\hat{P}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma}) &\leq 8 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, \sigma\} N} \left(1 + \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}}\right)} 1 \\ &\quad + 2 \sqrt{\frac{2 \log(\frac{18SAN}{\delta})}{\gamma^2 (1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} 1 + \frac{4 \log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N} 1 + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)^2} 1 \\ &\leq 10 \sqrt{\frac{2 \log(\frac{18SAN}{\delta})}{\gamma^3 (1-\gamma)^2 \max\{1-\gamma, \sigma\} N} \left(1 + \sqrt{\frac{\log(\frac{SAN}{\delta})}{(1-\gamma)^2 N}}\right)} 1 + \frac{5 \log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N} 1 \\ &\leq 160 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} 1 + \frac{5 \log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N} 1, \end{aligned} \quad (76)$$

where the last inequality holds by the fact  $\gamma \geq \frac{1}{4}$  and letting  $N \geq \frac{\log(\frac{SAN}{\delta})}{(1-\gamma)^2}$ .

**Step 3: controlling  $\|\hat{V}^{\pi^*, \sigma} - V^{\pi^*, \sigma}\|_{\infty}$ : bounding the second term in (66).** To proceed, applying Lemma 9 on the second term of the right hand side of (66) leads to

$$\begin{aligned} &(I - \gamma \hat{P}^{\pi^*, \hat{V}})^{-1} (\hat{P}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma}) \\ &\leq 2 (I - \gamma \hat{P}^{\pi^*, \hat{V}})^{-1} \left( \sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \sqrt{\text{Var}_{P^{\pi^*}}(V^{\star, \sigma})} + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} 1 \right) \\ &\leq \frac{2 \log(\frac{18SAN}{\delta})}{N(1-\gamma)} (I - \gamma \hat{P}^{\pi^*, \hat{V}})^{-1} 1 + 2 \underbrace{\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} (I - \gamma \hat{P}^{\pi^*, \hat{V}})^{-1} \sqrt{\text{Var}_{\hat{P}^{\pi^*, \hat{V}}}(\hat{V}^{\pi^*, \sigma})}}_{=: \mathcal{C}_4} \\ &\quad + 2 \underbrace{\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} (I - \gamma \hat{P}^{\pi^*, \hat{V}})^{-1} \left( \sqrt{\text{Var}_{\hat{P}^{\pi^*, \hat{V}}}(V^{\pi^*, \sigma} - \hat{V}^{\pi^*, \sigma})} \right)}_{=: \mathcal{C}_5} \\ &\quad + 2 \underbrace{\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} (I - \gamma \hat{P}^{\pi^*, \hat{V}})^{-1} \left( \sqrt{|\text{Var}_{\hat{P}^{\pi^*}}(V^{\star, \sigma}) - \text{Var}_{\hat{P}^{\pi^*, \hat{V}}}(V^{\star, \sigma})|} \right)}_{=: \mathcal{C}_6} \\ &\quad + 2 \underbrace{\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} (I - \gamma \hat{P}^{\pi^*, \hat{V}})^{-1} \left( \sqrt{\text{Var}_{P^{\pi^*}}(V^{\star, \sigma})} - \sqrt{\text{Var}_{\hat{P}^{\pi^*}}(V^{\star, \sigma})} \right)}_{=: \mathcal{C}_7}, \end{aligned} \quad (77)$$

where the last term  $\tilde{\mathcal{C}}_3$  can be controlled the same as  $\mathcal{C}_3$  in (75). We now bound the above terms separately.

- Applying Lemma 8 with  $P = \hat{P}^{\pi^*, \hat{V}}$ ,  $\pi = \pi^*$  and taking  $V = \hat{V}^{\pi^*, \sigma}$  which obeys  $\hat{V}^{\pi^*, \sigma} = r_{\pi^*} + \gamma \hat{P}^{\pi^*, \hat{V}} \hat{V}^{\pi^*, \sigma}$ , and in view of (70), the term  $\mathcal{C}_4$  in (77) can be controlled as follows:

$$\mathcal{C}_4 = 2 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} (I - \gamma \hat{P}^{\pi^*, \hat{V}})^{-1} \sqrt{\text{Var}_{\hat{P}^{\pi^*, \hat{V}}}(\hat{V}^{\pi^*, \sigma})}$$

$$\begin{aligned}
&\leq 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}}\sqrt{\frac{8(\max_s \widehat{V}^{\pi^*,\sigma}(s) - \min_s \widehat{V}^{\pi^*,\sigma}(s))}{\gamma^2(1-\gamma)^2}}1 \\
&\leq 8\sqrt{\frac{\log(\frac{18SAN}{\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}}1,
\end{aligned} \tag{78}$$

where the last inequality holds by applying Lemma 7.

- To continue, considering  $\mathcal{C}_5$ , we directly observe that (in view of (70))

$$\begin{aligned}
\mathcal{C}_5 &= 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}}\left(I - \gamma\widehat{P}^{\pi^*,\widehat{V}}\right)^{-1}\sqrt{\text{Var}_{\widehat{P}^{\pi^*,\widehat{V}}}(V^{\pi^*,\sigma} - \widehat{V}^{\pi^*,\sigma})} \\
&\leq 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2N}}\|V^{\pi^*,\sigma} - \widehat{V}^{\pi^*,\sigma}\|_{\infty}1.
\end{aligned} \tag{79}$$

- Then, it is easily verified that  $\mathcal{C}_6$  can be controlled similarly as (74) as follows:

$$\mathcal{C}_6 \leq 2\sqrt{\frac{2\log(\frac{18SAN}{\delta})}{\gamma^2(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}}1. \tag{80}$$

- Similarly,  $\mathcal{C}_7$  can be controlled the same as (75) shown below:

$$\mathcal{C}_7 \leq \frac{4\log(\frac{18SAN}{\delta})}{(1-\gamma)^2N}1. \tag{81}$$

Combining the results in (78), (79), (80), and (81) and inserting back to (77) leads to

$$\begin{aligned}
&\left(I - \gamma\widehat{P}^{\pi^*,\widehat{V}}\right)^{-1}\left(\widehat{P}^{\pi^*,V}V^{\pi^*,\sigma} - \underline{P}^{\pi^*,V}V^{\pi^*,\sigma}\right) \leq 8\sqrt{\frac{\log(\frac{18SAN}{\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}}1 \\
&\quad + 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2N}}\|V^{\pi^*,\sigma} - \widehat{V}^{\pi^*,\sigma}\|_{\infty}1 + 2\sqrt{\frac{2\log(\frac{18SAN}{\delta})}{\gamma^2(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}}1 + \frac{4\log(\frac{18SAN}{\delta})}{(1-\gamma)^2N}1 \\
&\leq 80\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}}1 + 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2N}}\|V^{\pi^*,\sigma} - \widehat{V}^{\pi^*,\sigma}\|_{\infty}1 + \frac{4\log(\frac{18SAN}{\delta})}{(1-\gamma)^2N}1,
\end{aligned} \tag{82}$$

where the last inequality follows from the assumption  $\gamma \geq \frac{1}{4}$ .

Finally, inserting (76) and (82) back to (66) yields

$$\begin{aligned}
\|\widehat{V}^{\pi^*,\sigma} - V^{\pi^*,\sigma}\|_{\infty} &\leq \max \left\{ 160\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} + \frac{5\log(\frac{18SAN}{\delta})}{(1-\gamma)^2N}, \right. \\
&\quad \left. 80\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} + 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2N}}\|V^{\pi^*,\sigma} - \widehat{V}^{\pi^*,\sigma}\|_{\infty} + \frac{4\log(\frac{18SAN}{\delta})}{(1-\gamma)^2N} \right\} \\
&\leq 160\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} + \frac{8\log(\frac{18SAN}{\delta})}{(1-\gamma)^2N},
\end{aligned} \tag{83}$$

where the last inequality holds by taking  $N \geq \frac{16\log(\frac{SAN}{\delta})}{(1-\gamma)^2}$ .



**Step 4: controlling  $\|\widehat{V}^{\widehat{\pi},\sigma} - V^{\widehat{\pi},\sigma}\|_\infty$ : bounding the first term in (67).** Unlike the earlier term, we now need to deal with the complicated statistical dependency between  $\widehat{\pi}$  and the empirical RMDP. To begin with, we introduce the following lemma which controls the main term on the right hand side of (67), which is proved in Appendix B.3.5.

**Lemma 12.** *Consider any  $\delta \in (0, 1)$ . Taking  $N \geq \log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)$ , with probability at least  $1 - \delta$ , one has*

$$\begin{aligned} \left| \widehat{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - \underline{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} \right| &\leq 2 \sqrt{\frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}^{\star,\sigma})} 1 + \frac{8 \log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{N(1-\gamma)} 1 + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} 1 \\ &\leq 10 \sqrt{\frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{(1-\gamma)^2 N}} 1 + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} 1. \end{aligned} \quad (84)$$

With Lemma 12 in hand, we have

$$\begin{aligned} &\left( I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \left( \widehat{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - \underline{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} \right) \\ &\stackrel{(i)}{\leq} \left( I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \left| \widehat{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - \underline{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} \right| \\ &\leq 2 \sqrt{\frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{N}} \left( I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \sqrt{\text{Var}_{P^{\widehat{\pi}}}(\widehat{V}^{\star,\sigma})} + \left( I - \gamma P_Q^{\widehat{\pi},V^{\widehat{\pi}}} \right)^{-1} \left( \frac{8 \log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{N(1-\gamma)} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \right) 1 \\ &\stackrel{(ii)}{\leq} \left( \frac{8 \log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{N(1-\gamma)^2} + \frac{2\gamma\varepsilon_{\text{opt}}}{(1-\gamma)^2} \right) 1 + \underbrace{2 \sqrt{\frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{N}} \left( I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\widehat{\pi},\sigma})}}_{=: \mathcal{D}_1} \\ &\quad + \underbrace{2 \sqrt{\frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{N}} \left( I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \sqrt{|\text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\star,\sigma}) - \text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\widehat{\pi},\sigma})|}}_{=: \mathcal{D}_2} \\ &\quad + \underbrace{2 \sqrt{\frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{N}} \left( I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \sqrt{|\text{Var}_{P^{\widehat{\pi}}}(\widehat{V}^{\star,\sigma}) - \text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\star,\sigma})|}}_{=: \mathcal{D}_3}, \end{aligned} \quad (85)$$

where (i) and (ii) hold by the fact that each row of  $(1-\gamma) \left( I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1}$  is a probability vector that falls into  $\Delta(\mathcal{S})$ .

The remainder of the proof will focus on controlling the three terms in (85) separately.

- For  $\mathcal{D}_1$ , we introduce the following lemma, whose proof is postponed to B.3.6.

**Lemma 13.** *Consider any  $\delta \in (0, 1)$ . Taking  $N \geq \frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{(1-\gamma)^2}$  and  $\varepsilon_{\text{opt}} \leq \frac{1-\gamma}{\gamma}$ , one has with probability at least  $1 - \delta$ ,*

$$\left( I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\widehat{\pi},\sigma})} \leq 6 \sqrt{\frac{1}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}}} 1 \leq 6 \sqrt{\frac{1}{(1-\gamma)^3 \gamma^2}} 1.$$

Applying Lemma 13 and (70) to (85) leads to

$$\begin{aligned} \mathcal{D}_1 &= 2 \sqrt{\frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{N}} \left( I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\widehat{\pi},\sigma})} \\ &\leq 12 \sqrt{\frac{\log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} 1. \end{aligned} \quad (86)$$

- Applying Lemma 2 with  $\|\widehat{V}^{\star,\sigma} - \widehat{V}^{\widehat{\pi},\sigma}\|_\infty \leq \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma}$  and (70),  $\mathcal{D}_2$  can be controlled as

$$\begin{aligned}\mathcal{D}_2 &= 2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}}\right)^{-1} \sqrt{\left|\text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\star,\sigma}) - \text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\widehat{\pi},\sigma})\right|} \\ &\leq 4\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}}\right)^{-1} \frac{\sqrt{\gamma\varepsilon_{\text{opt}}}}{1-\gamma} \leq 4\sqrt{\frac{\gamma\varepsilon_{\text{opt}} \log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^4 N}} 1.\end{aligned}\quad (87)$$

- $\mathcal{D}_3$  can be controlled similar to  $\mathcal{C}_2$  in (74) as follows:

$$\begin{aligned}\mathcal{D}_3 &= 2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}}\right)^{-1} \sqrt{\left|\text{Var}_{\underline{P}^{\widehat{\pi}}}(\widehat{V}^{\star,\sigma}) - \text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\star,\sigma})\right|} \\ &\leq 4\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}}\right)^{-1} \sqrt{\frac{1}{\gamma^2 \max\{1-\gamma, \sigma\}}} 1 \leq 4\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2(1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} 1\end{aligned}\quad (88)$$

Finally, summing up the results in (86), (87), and (88) and inserting them back to (85) yields: taking  $N \geq \frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2}$  and  $\varepsilon_{\text{opt}} \leq \frac{1-\gamma}{\gamma}$ , with probability at least  $1 - \delta$ ,

$$\begin{aligned}&\left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}}\right)^{-1} \left(\widehat{\underline{P}}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - \underline{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma}\right) \leq \left(\frac{8 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} + \frac{2\gamma\varepsilon_{\text{opt}}}{(1-\gamma)^2}\right) 1 \\ &\quad + 12\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} 1 + 4\sqrt{\frac{\gamma\varepsilon_{\text{opt}} \log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^4 N}} 1 + 4\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2(1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} 1 \\ &\leq 16\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\} N}} 1 + \frac{14 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} 1,\end{aligned}\quad (89)$$

where the last inequality holds by taking  $\varepsilon_{\text{opt}} \leq \min\left\{\frac{1-\gamma}{\gamma}, \frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma N}\right\} = \frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma N}$ .

**Step 5: controlling  $\|\widehat{V}^{\widehat{\pi},\sigma} - V^{\widehat{\pi},\sigma}\|_\infty$ : bounding the second term in (67).** Towards this, applying Lemma 12 leads to

$$\begin{aligned}&\left(I - \gamma \underline{P}^{\widehat{\pi},V}\right)^{-1} \left(\widehat{\underline{P}}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - \underline{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma}\right) \leq \left(I - \gamma \underline{P}^{\widehat{\pi},V}\right)^{-1} \left|\widehat{\underline{P}}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - \underline{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma}\right| \\ &\leq 2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\widehat{\pi},V}\right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\widehat{\pi}}}(\widehat{V}^{\star,\sigma})} + \left(I - \gamma \underline{P}^{\widehat{\pi},V}\right)^{-1} \left(\frac{8 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma}\right) 1 \\ &\leq \left(\frac{8 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} + \frac{2\gamma\varepsilon_{\text{opt}}}{(1-\gamma)^2}\right) 1 + \underbrace{2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\widehat{\pi},V}\right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\widehat{\pi},V}}(V^{\widehat{\pi},\sigma})}}_{=:\mathcal{D}_4} \\ &\quad + \underbrace{2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\widehat{\pi},V}\right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\widehat{\pi},V}}(\widehat{V}^{\widehat{\pi},\sigma} - V^{\widehat{\pi},\sigma})}}_{=:\mathcal{D}_5} \\ &\quad + \underbrace{2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left(I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}}\right)^{-1} \sqrt{\left|\text{Var}_{\underline{P}^{\widehat{\pi},V}}(\widehat{V}^{\star,\sigma}) - \text{Var}_{\underline{P}^{\widehat{\pi},V}}(\widehat{V}^{\widehat{\pi},\sigma})\right|}}_{=:\mathcal{D}_6}\end{aligned}$$

$$+ 2 \underbrace{\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left( I - \gamma \underline{P}^{\hat{\pi}, \hat{V}} \right)^{-1} \sqrt{\left| \text{Var}_{P^{\hat{\pi}}(\hat{V}^*, \sigma)} - \text{Var}_{\underline{P}^{\hat{\pi}, \sigma}}(\hat{V}^*, \sigma) \right|}}_{=: \mathcal{D}_7}. \quad (90)$$

We shall bound each of the terms separately.

- Applying Lemma 8 with  $P = \underline{P}^{\hat{\pi}, V}$ ,  $\pi = \hat{\pi}$ , and taking  $V = V^{\hat{\pi}, \sigma}$  which obeys  $V^{\hat{\pi}, \sigma} = r_{\hat{\pi}} + \gamma \underline{P}^{\hat{\pi}, V} V^{\hat{\pi}, \sigma}$ , the term  $\mathcal{D}_4$  can be controlled similar to (78) as follows:

$$\mathcal{D}_4 \leq 8 \sqrt{\frac{\log(\frac{54SAN^2}{\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1. \quad (91)$$

- For  $\mathcal{D}_5$ , it is observed that

$$\begin{aligned} \mathcal{D}_5 &= 2 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \left( I - \gamma \underline{P}^{\hat{\pi}, V} \right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi}, V}}(\hat{V}^{\hat{\pi}, \sigma} - V^{\hat{\pi}, \sigma})} \\ &\leq 2 \sqrt{\frac{\log(\frac{54SAN^2}{\delta})}{(1-\gamma)^2 N}} \left\| V^{\hat{\pi}, \sigma} - \hat{V}^{\hat{\pi}, \sigma} \right\|_{\infty} 1. \end{aligned} \quad (92)$$

- Next, observing that  $\mathcal{D}_6$  and  $\mathcal{D}_7$  are almost the same as the terms  $\mathcal{D}_2$  (controlled in (87)) and  $\mathcal{D}_3$  (controlled in (88)) in (85), it is easily verified that they can be controlled as follows

$$\mathcal{D}_6 \leq 4 \sqrt{\frac{\gamma \varepsilon_{\text{opt}} \log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^4 N}} 1, \quad \mathcal{D}_7 \leq 4 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1. \quad (93)$$

Then inserting the results in (91), (92), and (93) back to (90) leads to

$$\begin{aligned} &\left( I - \gamma \underline{P}^{\hat{\pi}, V} \right)^{-1} \left( \hat{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} \right) \leq \left( \frac{8 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} + \frac{2\gamma \varepsilon_{\text{opt}}}{(1-\gamma)^2} \right) 1 + 8 \sqrt{\frac{\log(\frac{54SAN^2}{\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1 \\ &\quad + 2 \sqrt{\frac{\log(\frac{54SAN^2}{\delta})}{(1-\gamma)^2 N}} \left\| V^{\hat{\pi}, \sigma} - \hat{V}^{\hat{\pi}, \sigma} \right\|_{\infty} 1 + 4 \sqrt{\frac{\gamma \varepsilon_{\text{opt}} \log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^4 N}} 1 + 4 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^2(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1 \\ &\leq 12 \sqrt{\frac{2 \log(\frac{8SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} 1 + \frac{14 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} 1 + 2 \sqrt{\frac{\log(\frac{54SAN^2}{\delta})}{(1-\gamma)^2 N}} \left\| V^{\hat{\pi}, \sigma} - \hat{V}^{\hat{\pi}, \sigma} \right\|_{\infty} 1, \end{aligned} \quad (94)$$

where the last inequality holds by letting  $\varepsilon_{\text{opt}} \leq \frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma N}$ , which directly satisfies  $\varepsilon_{\text{opt}} \leq \frac{1-\gamma}{\gamma}$  by letting  $N \geq \frac{\log(\frac{54SAN^2}{\delta})}{1-\gamma}$ .

Finally, inserting (89) and (94) back to (67) yields: taking  $\varepsilon_{\text{opt}} \leq \frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma N}$  and  $N \geq \frac{16 \log(\frac{54SAN^2}{\delta})}{(1-\gamma)^2}$ , with probability at least  $1 - \delta$ , one has

$$\begin{aligned} \left\| \hat{V}^{\hat{\pi}, \sigma} - V^{\hat{\pi}, \sigma} \right\|_{\infty} &\leq \max \left\{ 16 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} + \frac{14 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2}, \right. \\ &\quad \left. 12 \sqrt{\frac{2 \log(\frac{8SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} + \frac{14 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} + 2 \sqrt{\frac{\log(\frac{54SAN^2}{\delta})}{(1-\gamma)^2 N}} \left\| V^{\hat{\pi}, \sigma} - \hat{V}^{\hat{\pi}, \sigma} \right\|_{\infty} \right\} \\ &\leq 24 \sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} + \frac{28 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2}. \end{aligned} \quad (95)$$

**Step 6: summing up the results.** Summing up the results in (83) and (95) and inserting back to (61) complete the proof as follows: taking  $\varepsilon_{\text{opt}} \leq \frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma N}$  and  $N \geq \frac{16 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2}$ , with probability at least  $1 - \delta$ ,

$$\begin{aligned}
\|V^{\star, \sigma} - V^{\hat{\pi}, \sigma}\|_{\infty} &\leq \|V^{\pi^*, \sigma} - \hat{V}^{\pi^*, \sigma}\|_{\infty} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} + \|\hat{V}^{\hat{\pi}, \sigma} - V^{\hat{\pi}, \sigma}\|_{\infty} \\
&\leq \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} + 160\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} + \frac{8 \log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N} \\
&\quad + 24\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} + \frac{28 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} \\
&\leq 184\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}} + \frac{36 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)^2} \\
&\leq 1508\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2 \max\{1-\gamma, \sigma\}N}}, \tag{96}
\end{aligned}$$

where the last inequality holds by  $\gamma \geq \frac{1}{4}$  and  $N \geq \frac{16 \log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2}$ .

### B.3 Proof of the auxiliary lemmas

#### B.3.1 Proof of Lemma 7

To begin, note that there at leasts exist one state  $s_0$  for any  $V^{\pi, \sigma}$  such that  $V^{\pi, \sigma}(s_0) = \min_{s \in \mathcal{S}} V^{\pi, \sigma}(s)$ . With this in mind, for any policy  $\pi$ , one has by the definition in (5) and the Bellman's equation (7a),

$$\begin{aligned}
\max_{s \in \mathcal{S}} V^{\pi, \sigma}(s) &= \max_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi(\cdot | s)} \left[ r(s, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{s, a})} \mathcal{P} V^{\pi, \sigma} \right] \\
&\leq \max_{(s, a) \in \mathcal{S} \times \mathcal{A}} \left( 1 + \gamma \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{s, a})} \mathcal{P} V^{\pi, \sigma} \right),
\end{aligned}$$

where the second line holds since the reward function  $r(s, a) \in [0, 1]$  for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ . To continue, note that for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , there exists some  $\tilde{P}_{s, a} \in \mathbb{R}^{\mathcal{S}}$  constructed by reducing the values of some elements of  $P_{s, a}$  to obey  $P_{s, a} \geq \tilde{P}_{s, a} \geq 0$  and  $\sum_{s'} (P_{s, a}(s') - \tilde{P}_{s, a}(s')) = \sigma$ . This implies  $\tilde{P}_{s, a} + \sigma e_{s_0}^{\top} \in \mathcal{U}^{\sigma}(P_{s, a})$ , where  $e_{s_0}$  is the standard basis vector supported on  $s_0$ , since  $\frac{1}{2} \|\tilde{P}_{s, a} + \sigma e_{s_0}^{\top} - P_{s, a}\|_1 \leq \frac{1}{2} \|\tilde{P}_{s, a} - P_{s, a}\|_1 + \frac{\sigma}{2} = \sigma$ . Consequently,

$$\begin{aligned}
\inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{s, a})} \mathcal{P} V^{\pi, \sigma} &\leq \left( \tilde{P}_{s, a} + \sigma e_{s_0}^{\top} \right) V^{\pi, \sigma} \leq \|\tilde{P}_{s, a}\|_1 \|V^{\pi, \sigma}\|_{\infty} + \sigma V^{\pi, \sigma}(s_0) \\
&\leq (1 - \sigma) \max_{s \in \mathcal{S}} V^{\pi, \sigma}(s) + \sigma \min_{s \in \mathcal{S}} V^{\pi, \sigma}(s), \tag{97}
\end{aligned}$$

where the second inequality holds by  $\|\tilde{P}_{s, a}\|_1 = \sum_{s'} \tilde{P}_{s, a}(s') = -\sum_{s'} (P_{s, a}(s') - \tilde{P}_{s, a}(s')) + \sum_{s'} P_{s, a}(s') = 1 - \sigma$ . Plugging this back to the previous relation gives

$$\max_{s \in \mathcal{S}} V^{\pi, \sigma}(s) \leq 1 + \gamma(1 - \sigma) \max_{s \in \mathcal{S}} V^{\pi, \sigma}(s) + \gamma\sigma \min_{s \in \mathcal{S}} V^{\pi, \sigma}(s),$$

which, by rearranging terms, immediately yields

$$\begin{aligned}
\max_{s \in \mathcal{S}} V^{\pi, \sigma}(s) &\leq \frac{1 + \gamma\sigma \min_{s \in \mathcal{S}} V^{\pi, \sigma}(s)}{1 - \gamma(1 - \sigma)} \\
&\leq \frac{1}{(1 - \gamma) + \gamma\sigma} + \min_{s \in \mathcal{S}} V^{\pi, \sigma}(s) \leq \frac{1}{\gamma \max\{1 - \gamma, \sigma\}} + \min_{s \in \mathcal{S}} V^{\pi, \sigma}(s).
\end{aligned}$$

### B.3.2 Proof of Lemma 8

Observing that each row of  $P_\pi$  belongs to  $\Delta(S)$ , it can be directly verified that each row of  $(1-\gamma)(I - \gamma P_\pi)^{-1}$  falls into  $\Delta(S)$ . As a result,

$$\begin{aligned} (I - \gamma P_\pi)^{-1} \sqrt{\text{Var}_{P_\pi}(V^{\pi,P})} &= \frac{1}{1-\gamma} (1-\gamma) (I - \gamma P_\pi)^{-1} \sqrt{\text{Var}_{P_\pi}(V^{\pi,P})} \\ &\stackrel{(i)}{\leq} \frac{1}{1-\gamma} \sqrt{(1-\gamma) (I - \gamma P_\pi)^{-1} \text{Var}_{P_\pi}(V^{\pi,P})} \\ &= \sqrt{\frac{1}{1-\gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t (P_\pi)^t \text{Var}_{P_\pi}(V^{\pi,P})}, \end{aligned} \quad (98)$$

where (i) holds by Jensen's inequality.

To continue, denoting the minimum value of  $V$  as  $V_{\min} = \min_{s \in S} V^{\pi,P}(s)$  and  $V' := V^{\pi,P} - V_{\min}1$ . We control  $\text{Var}_{P_\pi}(V^{\pi,P})$  as follows:

$$\begin{aligned} \text{Var}_{P_\pi}(V^{\pi,P}) &\stackrel{(i)}{=} \text{Var}_{P_\pi}(V') = P_\pi(V' \circ V') - (P_\pi V') \circ (P_\pi V') \\ &\stackrel{(ii)}{=} P_\pi(V' \circ V') - \frac{1}{\gamma^2} (V' - r_\pi + (1-\gamma)V_{\min}1) \circ (V' - r_\pi + (1-\gamma)V_{\min}1) \\ &= P_\pi(V' \circ V') - \frac{1}{\gamma^2} V' \circ V' + \frac{2}{\gamma^2} V' \circ (r_\pi - (1-\gamma)V_{\min}1) - \frac{1}{\gamma^2} (r_\pi - (1-\gamma)V_{\min}1) \circ (r_\pi - (1-\gamma)V_{\min}1) \\ &\leq P_\pi(V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_\infty 1, \end{aligned} \quad (99)$$

where (i) holds by the fact that  $\text{Var}_{P_\pi}(V^{\pi,P} - b1) = \text{Var}_{P_\pi}(V^{\pi,P})$  for any scalar  $b$  and  $V^{\pi,P} \in \mathbb{R}^S$ , (ii) follows from  $V' = r_\pi + \gamma P_\pi V^{\pi,P} - V_{\min}1 = r_\pi - (1-\gamma)V_{\min}1 + \gamma P_\pi V'$ , and the last line arises from  $\frac{1}{\gamma^2} V' \circ V' \geq \frac{1}{\gamma} V' \circ V'$  and  $\|r_\pi - (1-\gamma)V_{\min}1\|_\infty \leq 1$ . Plugging (99) back to (98) leads to

$$\begin{aligned} (I - \gamma P_\pi)^{-1} \sqrt{\text{Var}_{P_\pi}(V^{\pi,P})} &\leq \sqrt{\frac{1}{1-\gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t (P_\pi)^t \left( P_\pi(V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_\infty 1 \right)} \\ &\stackrel{(i)}{\leq} \sqrt{\frac{1}{1-\gamma}} \sqrt{\left| \sum_{t=0}^{\infty} \gamma^t (P_\pi)^t \left( P_\pi(V' \circ V') - \frac{1}{\gamma} V' \circ V' \right) \right|} + \sqrt{\frac{1}{1-\gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t (P_\pi)^t \frac{2}{\gamma^2} \|V'\|_\infty 1} \\ &\leq \sqrt{\frac{1}{1-\gamma}} \sqrt{\left| \left( \sum_{t=0}^{\infty} \gamma^t (P_\pi)^{t+1} - \sum_{t=0}^{\infty} \gamma^{t-1} (P_\pi)^t \right) (V' \circ V') \right|} + \sqrt{\frac{2\|V'\|_\infty 1}{\gamma^2(1-\gamma)^2}} \\ &\stackrel{(ii)}{\leq} \sqrt{\frac{\|V'\|_\infty^2 1}{\gamma(1-\gamma)}} + \sqrt{\frac{2\|V'\|_\infty 1}{\gamma^2(1-\gamma)^2}} \\ &\leq \sqrt{\frac{8\|V'\|_\infty 1}{\gamma^2(1-\gamma)^2}}, \end{aligned} \quad (100)$$

where (i) holds by the triangle inequality, (ii) holds by following recursion, and the last inequality holds by  $\|V'\|_\infty \leq \frac{1}{1-\gamma}$ .

### B.3.3 Proof of Lemma 9

**Step 1: controlling the point-wise concentration.** We first consider a more general term w.r.t. any fixed (independent from  $\hat{P}^0$ ) value vector  $V$  obeying  $0 \leq V \leq \frac{1}{1-\gamma}1$  and any policy  $\pi$ . Invoking Lemma 4

leads to that for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,

$$\begin{aligned}
\left| \widehat{P}_{s,a}^{\pi,V} V - P_{s,a}^{\pi,V} V \right| &\leq \left| \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left\{ \widehat{P}_{s,a}^0 [V]_\alpha - \sigma \left( \alpha - \min_{s'} [V]_\alpha (s') \right) \right\} \right. \\
&\quad \left. - \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left\{ P_{s,a}^0 [V]_\alpha - w \sigma \left( \alpha - \min_{s'} [V]_\alpha (s') \right) \right\} \right| \\
&\leq \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \underbrace{\left| \left( P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [V]_\alpha \right|}_{=: g_{s,a}(\alpha, V)}, \tag{101}
\end{aligned}$$

where the last inequality holds by that the maximum operator is 1-Lipschitz.

Then for a fixed  $\alpha$  and any vector  $V$  that is independent with  $\widehat{P}^0$ , using the Bernstein's inequality, one has with probability at least  $1 - \delta$ ,

$$\begin{aligned}
g_{s,a}(\alpha, V) &= \left| \left( P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [V]_\alpha \right| \leq \sqrt{\frac{2 \log(\frac{2}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}([V]_\alpha)} + \frac{2 \log(\frac{2}{\delta})}{3N(1-\gamma)} \\
&\leq \sqrt{\frac{2 \log(\frac{2}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{2 \log(\frac{2}{\delta})}{3N(1-\gamma)}. \tag{102}
\end{aligned}$$

**Step 2: deriving the uniform concentration.** To obtain the union bound, we first notice that  $g_{s,a}(\alpha, V)$  is 1-Lipschitz w.r.t.  $\alpha$  for any  $V$  obeying  $\|V\|_\infty \leq \frac{1}{1-\gamma}$ . In addition, we can construct an  $\varepsilon_1$ -net  $N_{\varepsilon_1}$  over  $[0, \frac{1}{1-\gamma}]$  whose size satisfies  $|N_{\varepsilon_1}| \leq \frac{3}{\varepsilon_1(1-\gamma)}$  (Vershynin, 2018). By the union bound and (102), it holds with probability at least  $1 - \frac{\delta}{SA}$  that for all  $\alpha \in N_{\varepsilon_1}$ ,

$$g_{s,a}(\alpha, V) \leq \sqrt{\frac{2 \log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{2 \log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{3N(1-\gamma)}. \tag{103}$$

Combined with (101), it yields that,

$$\begin{aligned}
\left| \widehat{P}_{s,a}^{\pi,V} V - P_{s,a}^{\pi,V} V \right| &\leq \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left| \left( P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [V]_\alpha \right| \\
&\stackrel{(i)}{\leq} \varepsilon_1 + \sup_{\alpha \in N_{\varepsilon_1}} \left| \left( P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [V]_\alpha \right| \\
&\stackrel{(ii)}{\leq} \varepsilon_1 + \sqrt{\frac{2 \log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{2 \log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{3N(1-\gamma)} \tag{104}
\end{aligned}$$

$$\begin{aligned}
&\stackrel{(iii)}{\leq} \sqrt{\frac{2 \log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{\log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{N(1-\gamma)} \\
&\stackrel{(iv)}{\leq} 2 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V)} + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} \tag{105}
\end{aligned}$$

$$\begin{aligned}
&\leq 2 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \|V\|_\infty + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} \\
&\leq 3 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} \tag{106}
\end{aligned}$$

where (i) follows from that the optimal  $\alpha^*$  falls into the  $\varepsilon_1$ -ball centered around some point inside  $N_{\varepsilon_1}$  and  $g_{s,a}(\alpha, V)$  is 1-Lipschitz, (ii) holds by (103), (iii) arises from taking  $\varepsilon_1 = \frac{\log(\frac{2SA|N_{\varepsilon_1}|}{\delta})}{3N(1-\gamma)}$ , (iv) is verified by  $|N_{\varepsilon_1}| \leq \frac{3}{\varepsilon_1(1-\gamma)} \leq 9N$ , and the last inequality is due to the fact  $\|V^{*,\sigma}\|_\infty \leq \frac{1}{1-\gamma}$  and letting  $N \geq \log(\frac{18SAN}{\delta})$ .

To continue, applying (105) and (106) with  $\pi = \pi^*$  and  $V = V^{*,\sigma}$  (independent with  $\hat{P}^0$ ) and taking the union bound over  $(s, a) \in \mathcal{S} \times \mathcal{A}$  gives that with probability at least  $1 - \delta$ , it holds simultaneously for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$  that

$$\begin{aligned} \left| \hat{P}_{s,a}^{\pi^*,V} V^{*,\sigma} - P_{s,a}^{\pi^*,V} V^{*,\sigma} \right| &\leq 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(V^{*,\sigma})} + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} \\ &\leq 3\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}}. \end{aligned} \quad (107)$$

By converting (107) to the matrix form, one has with probability at least  $1 - \delta$ ,

$$\begin{aligned} \left| \underline{\hat{P}}^{\pi^*,V} V^{\pi^*,\sigma} - \underline{P}^{\pi^*,V} V^{\pi^*,\sigma} \right| &\leq 2\sqrt{\frac{\log(\frac{18SAN}{\delta})}{N}} \sqrt{\text{Var}_{P^{\pi^*}}(V^{*,\sigma})} + \frac{\log(\frac{18SAN}{\delta})}{N(1-\gamma)} 1 \\ &\leq 3\sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} 1. \end{aligned} \quad (108)$$

### B.3.4 Proof of Lemma 10

Following the same argument as (98), it follows

$$\left( I - \gamma \underline{\hat{P}}^{\pi^*,V} \right)^{-1} \sqrt{\text{Var}_{\underline{\hat{P}}^{\pi^*,V}}(V^{*,\sigma})} = \sqrt{\frac{1}{1-\gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t \left( \underline{\hat{P}}^{\pi^*,V} \right)^t \text{Var}_{\underline{\hat{P}}^{\pi^*,V}}(V^{*,\sigma})}. \quad (109)$$

To continue, we first focus on controlling  $\text{Var}_{\underline{\hat{P}}^{\pi^*,V}}(V^{*,\sigma})$ . Towards this, denoting the minimum value of  $V^{*,\sigma}$  as  $V_{\min} := \min_{s \in \mathcal{S}} V^{*,\sigma}(s)$  and  $V' := V^{*,\sigma} - V_{\min} 1$ , we arrive at (see the robust Bellman's consistency equation in (53))

$$\begin{aligned} V' &= V^{*,\sigma} - V_{\min} 1 = r_{\pi^*} + \gamma \underline{\hat{P}}^{\pi^*,V} V^{*,\sigma} - V_{\min} 1 \\ &= r_{\pi^*} + \gamma \underline{\hat{P}}^{\pi^*,V} V^{*,\sigma} + \gamma \left( \underline{P}^{\pi^*,V} - \underline{\hat{P}}^{\pi^*,V} \right) V^{*,\sigma} - V_{\min} 1 \\ &= r_{\pi^*} - (1-\gamma) V_{\min} 1 + \gamma \underline{\hat{P}}^{\pi^*,V} V' + \gamma \left( \underline{P}^{\pi^*,V} - \underline{\hat{P}}^{\pi^*,V} \right) V^{*,\sigma} \\ &= r'_{\pi^*} + \gamma \underline{\hat{P}}^{\pi^*,V} V' + \gamma \left( \underline{P}^{\pi^*,V} - \underline{\hat{P}}^{\pi^*,V} \right) V^{*,\sigma}, \end{aligned} \quad (110)$$

where the last line holds by letting  $r'_{\pi^*} := r_{\pi^*} - (1-\gamma) V_{\min} 1 \leq r_{\pi^*}$ . With the above fact in hand, we control  $\text{Var}_{\underline{\hat{P}}^{\pi^*,V}}(V^{*,\sigma})$  as follows:

$$\begin{aligned} \text{Var}_{\underline{\hat{P}}^{\pi^*,V}}(V^{*,\sigma}) &\stackrel{(i)}{=} \text{Var}_{\underline{\hat{P}}^{\pi^*,V}}(V') = \underline{\hat{P}}^{\pi^*,V} (V' \circ V') - (\underline{\hat{P}}^{\pi^*,V} V') \circ (\underline{\hat{P}}^{\pi^*,V} V') \\ &\stackrel{(ii)}{=} \underline{\hat{P}}^{\pi^*,V} (V' \circ V') - \frac{1}{\gamma^2} \left( V' - r'_{\pi^*} - \gamma \left( \underline{P}^{\pi^*,V} - \underline{\hat{P}}^{\pi^*,V} \right) V^{*,\sigma} \right)^{\circ 2} \\ &= \underline{\hat{P}}^{\pi^*,V} (V' \circ V') - \frac{1}{\gamma^2} V' \circ V' + \frac{2}{\gamma^2} V' \circ \left( r'_{\pi^*} + \gamma \left( \underline{P}^{\pi^*,V} - \underline{\hat{P}}^{\pi^*,V} \right) V^{*,\sigma} \right) \\ &\quad - \frac{1}{\gamma^2} \left( r'_{\pi^*} + \gamma \left( \underline{P}^{\pi^*,V} - \underline{\hat{P}}^{\pi^*,V} \right) V^{*,\sigma} \right)^{\circ 2} \\ &\stackrel{(iii)}{\leq} \underline{\hat{P}}^{\pi^*,V} (V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_{\infty} 1 + \frac{2}{\gamma} \|V'\|_{\infty} \left| \left( \underline{P}^{\pi^*,V} - \underline{\hat{P}}^{\pi^*,V} \right) V^{*,\sigma} \right| \end{aligned} \quad (111)$$

$$\leq \underline{\hat{P}}^{\pi^*,V} (V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_{\infty} 1 + \frac{6}{\gamma} \|V'\|_{\infty} \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1-\gamma)^2 N}} 1, \quad (112)$$



where (i) holds by the fact that  $\text{Var}_{P_\pi}(V - b1) = \text{Var}_{P_\pi}(V)$  for any scalar  $b$  and  $V \in \mathbb{R}^S$ , (ii) follows from (110), (iii) arises from  $\frac{1}{\gamma^2} V' \circ V' \geq \frac{1}{\gamma} V' \circ V'$  and  $-1 \leq r_{\pi^*} - (1 - \gamma)V_{\min}1 = r'_{\pi^*} \leq r_{\pi^*} \leq 1$ , and the last inequality holds by Lemma 9.

Plugging (112) into (109) leads to

$$\begin{aligned}
& \left( I - \gamma \hat{P}^{\pi^*, V} \right)^{-1} \sqrt{\text{Var}_{\hat{P}^{\pi^*, V}}(V^{\star, \sigma})} \\
& \leq \sqrt{\frac{1}{1 - \gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t \left( \hat{P}^{\pi^*, V} \right)^t \left( \hat{P}^{\pi^*, V} (V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_{\infty} 1 + \frac{6}{\gamma} \|V'\|_{\infty} \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1 - \gamma)^2 N}} 1 \right)} \\
& \stackrel{(i)}{\leq} \sqrt{\frac{1}{1 - \gamma}} \sqrt{\left| \sum_{t=0}^{\infty} \gamma^t \left( \hat{P}^{\pi^*, V} \right)^t \left( \hat{P}^{\pi^*, V} (V' \circ V') - \frac{1}{\gamma} V' \circ V' \right) \right|} \\
& \quad + \sqrt{\frac{1}{1 - \gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t \left( \hat{P}^{\pi^*, V} \right)^t \left( \frac{2}{\gamma^2} \|V'\|_{\infty} 1 + \frac{6}{\gamma} \|V'\|_{\infty} \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1 - \gamma)^2 N}} 1 \right)} \\
& \leq \sqrt{\frac{1}{1 - \gamma}} \sqrt{\left| \sum_{t=0}^{\infty} \gamma^t \left( \hat{P}^{\pi^*, V} \right)^t \left[ \hat{P}^{\pi^*, V} (V' \circ V') - \frac{1}{\gamma} V' \circ V' \right] \right|} + \sqrt{\frac{\left( 2 + 6 \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1 - \gamma)^2 N}} \right) \|V'\|_{\infty}}{(1 - \gamma)^2 \gamma^2}} 1, \quad (113)
\end{aligned}$$

where (i) holds by the triangle inequality. Therefore, the remainder of the proof shall focus on the first term, which follows

$$\begin{aligned}
& \left| \sum_{t=0}^{\infty} \gamma^t \left( \hat{P}^{\pi^*, V} \right)^t \left( \hat{P}^{\pi^*, V} (V' \circ V') - \frac{1}{\gamma} V' \circ V' \right) \right| \\
& = \left| \left( \sum_{t=0}^{\infty} \gamma^t \left( \hat{P}^{\pi^*, V} \right)^{t+1} - \sum_{t=0}^{\infty} \gamma^{t-1} \left( \hat{P}^{\pi^*, V} \right)^t \right) (V' \circ V') \right| \leq \frac{1}{\gamma} \|V'\|_{\infty}^2 1 \quad (114)
\end{aligned}$$

by recursion. Inserting (114) back to (113) leads to

$$\begin{aligned}
& \left( I - \gamma \hat{P}^{\pi^*, V} \right)^{-1} \sqrt{\text{Var}_{\hat{P}^{\pi^*, V}}(V^{\star, \sigma})} \\
& \leq \sqrt{\frac{\|V'\|_{\infty}^2}{\gamma(1 - \gamma)}} 1 + 3 \sqrt{\frac{\left( 1 + \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1 - \gamma)^2 N}} \right) \|V'\|_{\infty}}{(1 - \gamma)^2 \gamma^2}} 1 \\
& \leq 4 \sqrt{\frac{\left( 1 + \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1 - \gamma)^2 N}} \right) \|V'\|_{\infty}}{(1 - \gamma)^2 \gamma^2}} 1 \leq 4 \sqrt{\frac{\left( 1 + \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1 - \gamma)^2 N}} \right)}{\gamma^3 (1 - \gamma)^2 \max\{1 - \gamma, \sigma\}}} 1 \leq 4 \sqrt{\frac{\left( 1 + \sqrt{\frac{\log(\frac{18SAN}{\delta})}{(1 - \gamma)^2 N}} \right)}{\gamma^3 (1 - \gamma)^3}} 1, \quad (115)
\end{aligned}$$

where the penultimate inequality follows from applying Lemma 7 with  $P = P^0$  and  $\pi = \pi^*$ :

$$\|V'\|_{\infty} = \max_{s \in \mathcal{S}} V^{\star, \sigma}(s) - \min_{s \in \mathcal{S}} V^{\star, \sigma}(s) \leq \frac{1}{\gamma \max\{1 - \gamma, \sigma\}}.$$

### B.3.5 Proof of Lemma 12

To begin with, for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , invoking the results in (101), we have

$$\begin{aligned}
& \left| \hat{P}_{s,a}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - P_{s,a}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} \right| \leq \max_{\alpha \in [\min_s \hat{V}^{\hat{\pi}, \sigma}(s), \max_s \hat{V}^{\hat{\pi}, \sigma}(s)]} \left| \left( P_{s,a}^0 - \hat{P}_{s,a}^0 \right) [\hat{V}^{\hat{\pi}, \sigma}]_{\alpha} \right| \\
& \stackrel{(i)}{\leq} \max_{\alpha \in [\min_s \hat{V}^{\hat{\pi}, \sigma}(s), \max_s \hat{V}^{\hat{\pi}, \sigma}(s)]} \left( \left| \left( P_{s,a}^0 - \hat{P}_{s,a}^0 \right) [\hat{V}^{\star, \sigma}]_{\alpha} \right| + \left| \left( P_{s,a}^0 - \hat{P}_{s,a}^0 \right) \left( [\hat{V}^{\hat{\pi}, \sigma}]_{\alpha} - [\hat{V}^{\star, \sigma}]_{\alpha} \right) \right| \right)
\end{aligned}$$

$$\begin{aligned}
&\leq \max_{\alpha \in [\min_s \hat{V}^{\hat{\pi}, \sigma}(s), \max_s \hat{V}^{\hat{\pi}, \sigma}(s)]} \left( \left| \left( P_{s,a}^0 - \hat{P}_{s,a}^0 \right) [\hat{V}^{\star, \sigma}]_{\alpha} \right| + \left\| P_{s,a}^0 - \hat{P}_{s,a}^0 \right\|_1 \left\| [\hat{V}^{\hat{\pi}, \sigma}]_{\alpha} - [\hat{V}^{\star, \sigma}]_{\alpha} \right\|_{\infty} \right) \\
&\stackrel{(ii)}{\leq} \max_{\alpha \in [\min_s \hat{V}^{\hat{\pi}, \sigma}(s), \max_s \hat{V}^{\hat{\pi}, \sigma}(s)]} \left| \left( P_{s,a}^0 - \hat{P}_{s,a}^0 \right) [\hat{V}^{\star, \sigma}]_{\alpha} \right| + 2 \left\| \hat{V}^{\hat{\pi}, \sigma} - \hat{V}^{\star, \sigma} \right\|_{\infty} \\
&\stackrel{(iii)}{\leq} \max_{\alpha \in [\min_s \hat{V}^{\hat{\pi}, \sigma}(s), \max_s \hat{V}^{\hat{\pi}, \sigma}(s)]} \left| \left( P_{s,a}^0 - \hat{P}_{s,a}^0 \right) [\hat{V}^{\star, \sigma}]_{\alpha} \right| + \frac{2\gamma\epsilon_{\text{opt}}}{1-\gamma}, \tag{116}
\end{aligned}$$

where (i) holds by the triangle inequality, and (ii) follows from  $\|P_{s,a}^0 - \hat{P}_{s,a}^0\|_1 \leq 2$  and  $\|[\hat{V}^{\hat{\pi}, \sigma}]_{\alpha} - [\hat{V}^{\star, \sigma}]_{\alpha}\|_{\infty} \leq \|\hat{V}^{\hat{\pi}, \sigma} - \hat{V}^{\star, \sigma}\|_{\infty}$ , and (iii) follows from (60).

To control  $\left| \left( P_{s,a}^0 - \hat{P}_{s,a}^0 \right) [\hat{V}^{\star, \sigma}]_{\alpha} \right|$  in (116) for any given  $\alpha \in [0, \frac{1}{1-\gamma}]$ , and tame the dependency between  $\hat{V}^{\star, \sigma}$  and  $\hat{P}^0$ , we resort to the following leave-one-out argument motivated by (Agarwal et al., 2020; Li et al., 2022b; Shi and Chi, 2022). Specifically, we first construct a set of auxiliary RMDPs which simultaneously have the desired statistical independence between robust value functions and the estimated nominal transition kernel, and are minimally different from the original RMDPs under consideration. Then we control the term of interest associated with these auxiliary RMDPs and show the value is close to the target quantity for the desired RMDP. The process is divided into several steps as below.

**Step 1: construction of auxiliary RMDPs with deterministic empirical nominal transitions.**

Recall that we target the empirical infinite-horizon robust MDP  $\widehat{\mathcal{M}}_{\text{rob}}$  with the nominal transition kernel  $\hat{P}^0$ . Towards this, we can construct an auxiliary robust MDP  $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$  for each state  $s$  and any non-negative scalar  $u \geq 0$ , so that it is the same as  $\widehat{\mathcal{M}}_{\text{rob}}$  except for the transition properties in state  $s$ . In particular, we define the nominal transition kernel and reward function of  $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$  as  $P^{s,u}$  and  $r^{s,u}$ , which are expressed as follows

$$\begin{cases} P^{s,u}(s' | s, a) = \mathbb{1}(s' = s) & \text{for all } (s', a) \in \mathcal{S} \times \mathcal{A}, \\ P^{s,u}(\cdot | \tilde{s}, a) = \hat{P}^0(\cdot | \tilde{s}, a) & \text{for all } (\tilde{s}, a) \in \mathcal{S} \times \mathcal{A} \text{ and } \tilde{s} \neq s, \end{cases} \tag{117}$$

and

$$\begin{cases} r^{s,u}(s, a) = u & \text{for all } a \in \mathcal{A}, \\ r^{s,u}(\tilde{s}, a) = r(\tilde{s}, a) & \text{for all } (\tilde{s}, a) \in \mathcal{S} \times \mathcal{A} \text{ and } \tilde{s} \neq s. \end{cases} \tag{118}$$

It is evident that the nominal transition probability at state  $s$  of the auxiliary  $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ , i.e. it never leaves state  $s$  once entered. This useful property removes the randomness of  $\hat{P}_{s,a}^0$  for all  $a \in \mathcal{A}$  in state  $s$ , which will be leveraged later.

Correspondingly, the robust Bellman operator  $\hat{\mathcal{T}}_{s,u}^{\sigma}(\cdot)$  associated with the RMDP  $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$  is defined as

$$\forall (\tilde{s}, a) \in \mathcal{S} \times \mathcal{A}: \quad \hat{\mathcal{T}}_{s,u}^{\sigma}(Q)(\tilde{s}, a) = r^{s,u}(\tilde{s}, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{\tilde{s},a}^{s,u})} \mathcal{P}V, \quad \text{with } V(\tilde{s}) = \max_a Q(\tilde{s}, a). \tag{119}$$

**Step 2: fixed-point equivalence between  $\widehat{\mathcal{M}}_{\text{rob}}$  and the auxiliary RMDP  $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$ .** Recall that  $\hat{Q}^{\star, \sigma}$  is the unique fixed point of  $\hat{\mathcal{T}}^{\sigma}(\cdot)$  with the corresponding robust value  $\hat{V}^{\star, \sigma}$ . We assert that the corresponding robust value function  $\hat{V}_{s,u}^{\star, \sigma}$  obtained from the fixed point of  $\hat{\mathcal{T}}_{s,u}^{\sigma}(\cdot)$  aligns with the robust value function  $\hat{V}^{\star, \sigma}$  derived from  $\hat{\mathcal{T}}^{\sigma}(\cdot)$ , as long as we choose  $u$  in the following manner:

$$u^{\star} := u^{\star}(s) = \hat{V}^{\star, \sigma}(s) - \gamma \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(e_s)} \mathcal{P}\hat{V}^{\star, \sigma}. \tag{120}$$

where  $e_s$  is the  $s$ -th standard basis vector in  $\mathbb{R}^{\mathcal{S}}$ . Towards verifying this, we shall break our arguments in two different cases.

- **For state  $s$ :** One has for any  $a \in \mathcal{A}$ :

$$\begin{aligned} r^{s,u^*}(s,a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s,a}^{s,u^*})} \mathcal{P}\widehat{V}^{*,\sigma} &= u^* + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(e_s)} \mathcal{P}\widehat{V}^{*,\sigma} \\ &= \widehat{V}^{*,\sigma}(s) - \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(e_s)} \mathcal{P}\widehat{V}^{*,\sigma} + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(e_s)} \mathcal{P}\widehat{V}^{*,\sigma} = \widehat{V}^{*,\sigma}(s), \end{aligned} \quad (121)$$

where the first equality follows from the definition of  $P_{s,a}^{s,u^*}$  in (117), and the second equality follows from plugging in the definition of  $u^*$  in (120).

- **For state  $s' \neq s$ :** It is easily verified that for all  $a \in \mathcal{A}$ ,

$$\begin{aligned} r^{s,u^*}(s',a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{s',a}^{s,u^*})} \mathcal{P}\widehat{V}^{*,\sigma} &= r(s',a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(\widehat{P}_{s',a}^0)} \mathcal{P}\widehat{V}^{*,\sigma} \\ &= \widehat{\mathcal{T}}^\sigma(\widehat{Q}^{*,\sigma})(s',a) = \widehat{Q}^{*,\sigma}(s',a), \end{aligned} \quad (122)$$

where the first equality follows from the definitions in (118) and (117), and the last line arises from the definition of the robust Bellman operator in (15), and that  $\widehat{Q}^{*,\sigma}$  is the fixed point of  $\widehat{\mathcal{T}}^\sigma(\cdot)$  (see Lemma 3).

Combining the facts in the above two cases, we establish that there exists a fixed point  $\widehat{Q}_{s,u^*}^{*,\sigma}$  of the operator  $\widehat{\mathcal{T}}_{s,u^*}^\sigma(\cdot)$  by taking

$$\begin{cases} \widehat{Q}_{s,u^*}^{*,\sigma}(s,a) = \widehat{V}^{*,\sigma}(s) & \text{for all } a \in \mathcal{A}, \\ \widehat{Q}_{s,u^*}^{*,\sigma}(s',a) = \widehat{Q}^{*,\sigma}(s',a) & \text{for all } s' \neq s \text{ and } a \in \mathcal{A}. \end{cases} \quad (123)$$

Consequently, we confirm the existence of a fixed point of the operator  $\widehat{\mathcal{T}}_{s,u^*}^\sigma(\cdot)$ . In addition, its corresponding value function  $\widehat{V}_{s,u^*}^{*,\sigma}$  also coincides with  $\widehat{V}^{*,\sigma}$ . Note that the corresponding facts between  $\widehat{\mathcal{M}}_{\text{rob}}$  and  $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$  in Step 1 and step 2 holds in fact for any uncertainty set.

**Step 3: building an  $\varepsilon$ -net for all reward values  $u$ .** It is easily verified that

$$0 \leq u^* \leq \widehat{V}^{*,\sigma}(s) \leq \frac{1}{1-\gamma}. \quad (124)$$

We can construct a  $N_{\varepsilon_2}$ -net over the interval  $[0, \frac{1}{1-\gamma}]$ , where the size is bounded by  $|N_{\varepsilon_2}| \leq \frac{3}{\varepsilon_2(1-\gamma)}$  (Ver-shynin, 2018). Following the same arguments in the proof of Lemma 3, we can demonstrate that for each  $u \in N_{\varepsilon_2}$ , there exists a unique fixed point  $\widehat{Q}_{s,u}^{*,\sigma}$  of the operator  $\widehat{\mathcal{T}}_{s,u}^\sigma(\cdot)$ , which satisfies  $0 \leq \widehat{Q}_{s,u}^{*,\sigma} \leq \frac{1}{1-\gamma} \cdot 1$ .

Consequently, the corresponding robust value function also satisfies  $\left\| \widehat{V}_{s,u}^{*,\sigma} \right\|_\infty \leq \frac{1}{1-\gamma}$ .

By the definitions in (117) and (118), we observe that for all  $u \in N_{\varepsilon_2}$ ,  $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$  is statistically independent from  $\widehat{P}_{s,a}^0$ . This independence indicates that  $[\widehat{V}_{s,u}^{*,\sigma}]_\alpha$  and  $\widehat{P}_{s,a}^0$  are independent for a fixed  $\alpha$ . With this in mind, invoking the fact in (105) and (106) and taking the union bound over all  $(s,a,\alpha) \in \mathcal{S} \times \mathcal{A} \times N_{\varepsilon_1}$ ,  $u \in N_{\varepsilon_2}$  yields that, with probability at least  $1 - \delta$ , it holds for all  $(s,a,u) \in \mathcal{S} \times \mathcal{A} \times N_{\varepsilon_2}$  that

$$\begin{aligned} \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \left( P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}_{s,u}^{*,\sigma}]_\alpha \right| &\leq \varepsilon_2 + 2 \sqrt{\frac{\log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}_{s,u}^{*,\sigma})} + \frac{2 \log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{3N(1-\gamma)} \\ &\leq \varepsilon_2 + 3 \sqrt{\frac{\log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{(1-\gamma)^2 N}}, \end{aligned} \quad (125)$$

where the last inequality holds by the fact  $\text{Var}_{P_{s,a}^0}(\widehat{V}_{s,u}^{*,\sigma}) \leq \|\widehat{V}_{s,u}^{*,\sigma}\|_\infty \leq \frac{1}{1-\gamma}$  and letting  $N \geq \log\left(\frac{18SAN|N_{\varepsilon_2}|}{\delta}\right)$ .

**Step 4: uniform concentration.** Recalling that  $u^* \in [0, \frac{1}{1-\gamma}]$  (see (124)), we can always find some  $\bar{u} \in N_{\varepsilon_2}$  such that  $|\bar{u} - u^*| \leq \varepsilon_2$ . Consequently, plugging in the operator  $\hat{\mathcal{T}}_{s,u}^\sigma(\cdot)$  in (119) yields

$$\forall Q \in \mathbb{R}^{SA} : \quad \left\| \hat{\mathcal{T}}_{s,\bar{u}}^\sigma(Q) - \hat{\mathcal{T}}_{s,u^*}^\sigma(Q) \right\|_\infty = |\bar{u} - u^*| \leq \varepsilon_2$$

With this in mind, we observe that the fixed points of  $\hat{\mathcal{T}}_{s,\bar{u}}^\sigma(\cdot)$  and  $\hat{\mathcal{T}}_{s,u^*}^\sigma(\cdot)$  obey

$$\begin{aligned} \left\| \hat{Q}_{s,\bar{u}}^{\star,\sigma} - \hat{Q}_{s,u^*}^{\star,\sigma} \right\|_\infty &= \left\| \hat{\mathcal{T}}_{s,\bar{u}}^\sigma(\hat{Q}_{s,\bar{u}}^{\star,\sigma}) - \hat{\mathcal{T}}_{s,u^*}^\sigma(\hat{Q}_{s,u^*}^{\star,\sigma}) \right\|_\infty \\ &\leq \left\| \hat{\mathcal{T}}_{s,\bar{u}}^\sigma(\hat{Q}_{s,\bar{u}}^{\star,\sigma}) - \hat{\mathcal{T}}_{s,\bar{u}}^\sigma(\hat{Q}_{s,u^*}^{\star,\sigma}) \right\|_\infty + \left\| \hat{\mathcal{T}}_{s,\bar{u}}^\sigma(\hat{Q}_{s,u^*}^{\star,\sigma}) - \hat{\mathcal{T}}_{s,u^*}^\sigma(\hat{Q}_{s,u^*}^{\star,\sigma}) \right\|_\infty \\ &\leq \gamma \left\| \hat{Q}_{s,\bar{u}}^{\star,\sigma} - \hat{Q}_{s,u^*}^{\star,\sigma} \right\|_\infty + \varepsilon_2, \end{aligned}$$

where the last inequality holds by the fact that  $\hat{\mathcal{T}}_{s,u}^\sigma(\cdot)$  is a  $\gamma$ -contraction. It directly indicates that

$$\left\| \hat{Q}_{s,\bar{u}}^{\star,\sigma} - \hat{Q}_{s,u^*}^{\star,\sigma} \right\|_\infty \leq \frac{\varepsilon_2}{(1-\gamma)} \quad \text{and} \quad \left\| \hat{V}_{s,\bar{u}}^{\star,\sigma} - \hat{V}_{s,u^*}^{\star,\sigma} \right\|_\infty \leq \left\| \hat{Q}_{s,\bar{u}}^{\star,\sigma} - \hat{Q}_{s,u^*}^{\star,\sigma} \right\|_\infty \leq \frac{\varepsilon_2}{(1-\gamma)}. \quad (126)$$

Armed with the above facts, to control the first term in (116), invoking the identity  $\hat{V}^{\star,\sigma} = \hat{V}_{s,u^*}^{\star,\sigma}$  established in Step 2 gives that: for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,

$$\begin{aligned} &\max_{\alpha \in [\min_s \hat{V}^{\pi,\sigma}(s), \max_s \hat{V}^{\pi,\sigma}(s)]} \left| \left( P_{s,a}^0 - \hat{P}_{s,a}^0 \right) [\hat{V}^{\star,\sigma}]_\alpha \right| \\ &\leq \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \left( P_{s,a}^0 - \hat{P}_{s,a}^0 \right) [\hat{V}^{\star,\sigma}]_\alpha \right| = \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \left( P_{s,a}^0 - \hat{P}_{s,a}^0 \right) [\hat{V}_{s,u^*}^{\star,\sigma}]_\alpha \right| \\ &\stackrel{(i)}{\leq} \max_{\alpha \in [0, 1/(1-\gamma)]} \left\{ \left| \left( P_{s,a}^0 - \hat{P}_{s,a}^0 \right) [\hat{V}_{s,\bar{u}}^{\star,\sigma}]_\alpha \right| + \left| \left( P_{s,a}^0 - \hat{P}_{s,a}^0 \right) \left( [\hat{V}_{s,\bar{u}}^{\star,\sigma}]_\alpha - [\hat{V}_{s,u^*}^{\star,\sigma}]_\alpha \right) \right| \right\} \\ &\stackrel{(ii)}{\leq} \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \left( P_{s,a}^0 - \hat{P}_{s,a}^0 \right) [\hat{V}_{s,\bar{u}}^{\star,\sigma}]_\alpha \right| + \frac{2\varepsilon_2}{(1-\gamma)} \\ &\stackrel{(iii)}{\leq} \frac{2\varepsilon_2}{(1-\gamma)} + \varepsilon_2 + 2\sqrt{\frac{\log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\hat{V}_{s,\bar{u}}^{\star,\sigma})} + \frac{2\log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{3N(1-\gamma)} \\ &\leq \frac{3\varepsilon_2}{(1-\gamma)} + 2\sqrt{\frac{\log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\hat{V}^{\star,\sigma})} + \frac{2\log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{3N(1-\gamma)} \\ &\quad + 2\sqrt{\frac{\log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{N}} \sqrt{\left| \text{Var}_{P_{s,a}^0}(\hat{V}^{\star,\sigma}) - \text{Var}_{P_{s,a}^0}(\hat{V}_{s,\bar{u}}^{\star,\sigma}) \right|} \\ &\stackrel{(iv)}{\leq} \frac{3\varepsilon_2}{(1-\gamma)} + 2\sqrt{\frac{\log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\hat{V}^{\star,\sigma})} + \frac{2\log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{3N(1-\gamma)} + 2\sqrt{\frac{2\varepsilon_2 \log(\frac{18SAN|N_{\varepsilon_2}|}{\delta})}{N(1-\gamma)^2}} \\ &\leq 2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\hat{V}^{\star,\sigma})} + \frac{8\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)} \end{aligned} \quad (127)$$

$$\leq 10\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2N}}, \quad (128)$$

where (i) holds by the triangle inequality, (ii) arises from (the last inequality holds by (126))

$$\begin{aligned} \left| \left( P_{s,a}^0 - \hat{P}_{s,a}^0 \right) \left( [\hat{V}_{s,\bar{u}}^{\star,\sigma}]_\alpha - [\hat{V}_{s,u^*}^{\star,\sigma}]_\alpha \right) \right| &\leq \left\| P_{s,a}^0 - \hat{P}_{s,a}^0 \right\|_1 \left\| [\hat{V}_{s,\bar{u}}^{\star,\sigma}]_\alpha - [\hat{V}_{s,u^*}^{\star,\sigma}]_\alpha \right\|_\infty \\ &\leq 2 \left\| \hat{V}_{s,\bar{u}}^{\star,\sigma} - \hat{V}_{s,u^*}^{\star,\sigma} \right\|_\infty \leq \frac{2\varepsilon_2}{(1-\gamma)}, \end{aligned} \quad (129)$$

(iii) follows from (125), (iv) can be verified by applying Lemma 2 with (126). Here, the penultimate inequality holds by letting  $\varepsilon_2 = \frac{\log(\frac{18SAN|N\varepsilon_2|}{N})}{N}$ , which leads to  $|N_{\varepsilon_2}| \leq \frac{3}{\varepsilon_2(1-\gamma)} \leq \frac{3N}{1-\gamma}$ , and the last inequality holds by the fact  $\text{Var}_{P_{s,a}^0}(\widehat{V}^{\star,\sigma}) \leq \|\widehat{V}^{\star,\sigma}\|_\infty \leq \frac{1}{1-\gamma}$  and letting  $N \geq \log\left(\frac{54SAN^2}{(1-\gamma)\delta}\right)$ .

**Step 5: finishing up.** Inserting (127) and (128) back into (116) and combining with (128) give that with probability at least  $1 - \delta$ ,

$$\begin{aligned} \left| \widehat{P}_{s,a}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - P_{s,a}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} \right| &\leq \max_{\alpha \in [\min_s \widehat{V}^{\widehat{\pi},\sigma}(s), \max_s \widehat{V}^{\widehat{\pi},\sigma}(s)]} \left| \left( P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{\star,\sigma}]_\alpha \right| + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \\ &\leq \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \left( P_{s,a}^0 - \widehat{P}_{s,a}^0 \right) [\widehat{V}^{\star,\sigma}]_\alpha \right| + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \\ &\leq 2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}^{\star,\sigma})} + \frac{8\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \\ &\leq 10\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2N}} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \end{aligned} \quad (130)$$

holds for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ .

Finally, we complete the proof by compiling everything into the matrix form as follows:

$$\begin{aligned} \left| \widehat{\underline{P}}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} - \underline{P}^{\widehat{\pi},\widehat{V}} \widehat{V}^{\widehat{\pi},\sigma} \right| &\leq 2\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N}} \sqrt{\text{Var}_{P_{s,a}^0}(\widehat{V}^{\star,\sigma})} 1 + \frac{8\log(\frac{54SAN^2}{(1-\gamma)\delta})}{N(1-\gamma)} 1 + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} 1 \\ &\leq 10\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2N}} 1 + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} 1. \end{aligned} \quad (131)$$

### B.3.6 Proof of Lemma 13

The proof can be achieved by directly applying the same routine as Appendix B.3.4. Towards this, similar to (109), we arrive at

$$\left( I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\widehat{\pi},\sigma})} \leq \sqrt{\frac{1}{1-\gamma}} \sqrt{\sum_{t=0}^{\infty} \gamma^t \left( \underline{P}^{\widehat{\pi},\widehat{V}} \right)^t \text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\widehat{\pi},\sigma})}. \quad (132)$$

To control  $\text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\widehat{\pi},\sigma})$ , we denote the minimum value of  $\widehat{V}^{\widehat{\pi},\sigma}$  as  $V_{\min} = \min_{s \in \mathcal{S}} \widehat{V}^{\widehat{\pi},\sigma}(s)$  and  $V' := \widehat{V}^{\widehat{\pi},\sigma} - V_{\min} 1$ . By the same argument as (111), we arrive at

$$\begin{aligned} \text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\widehat{\pi},\sigma}) &\leq \underline{P}^{\widehat{\pi},\widehat{V}} (V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_\infty 1 + \frac{2}{\gamma} \|V'\|_\infty \left| \left( \widehat{\underline{P}}^{\widehat{\pi},\widehat{V}} - \underline{P}^{\widehat{\pi},\widehat{V}} \right) \widehat{V}^{\widehat{\pi},\sigma} \right| \\ &\leq \underline{P}^{\widehat{\pi},\widehat{V}} (V' \circ V') - \frac{1}{\gamma} V' \circ V' + \frac{2}{\gamma^2} \|V'\|_\infty 1 + \frac{2}{\gamma} \|V'\|_\infty \left( 10\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2N}} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \right) 1, \end{aligned} \quad (133)$$

where the last inequality makes use of Lemma 12. Plugging (133) back into (132) leads to

$$\begin{aligned} \left( I - \gamma \underline{P}^{\widehat{\pi},\widehat{V}} \right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\widehat{\pi},\widehat{V}}}(\widehat{V}^{\widehat{\pi},\sigma})} &\stackrel{(i)}{\leq} \sqrt{\frac{1}{1-\gamma}} \sqrt{\left| \sum_{t=0}^{\infty} \gamma^t \left( \underline{P}^{\widehat{\pi},\widehat{V}} \right)^t \left( \underline{P}^{\widehat{\pi},\widehat{V}} (V' \circ V') - \frac{1}{\gamma} V' \circ V' \right) \right|} \\ &\quad + \sqrt{\frac{1}{(1-\gamma)^2\gamma^2} \left( 2 + 20\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)\delta})}{(1-\gamma)^2N}} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} \right) \|V'\|_\infty 1} \end{aligned}$$

$$\begin{aligned}
&\stackrel{(ii)}{\leq} \sqrt{\frac{\|V'\|_\infty^2}{\gamma(1-\gamma)}} 1 + \sqrt{\frac{\left(2 + 20\sqrt{\frac{\log(\frac{54SAN^2}{(1-\gamma)^\delta})}{(1-\gamma)^2N}} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma}\right)\|V'\|_\infty}{(1-\gamma)^2\gamma^2}} 1 \\
&\stackrel{(iii)}{\leq} \sqrt{\frac{\|V'\|_\infty^2}{\gamma(1-\gamma)}} 1 + \sqrt{\frac{24\|V'\|_\infty}{(1-\gamma)^2\gamma^2}} 1 \leq 6\sqrt{\frac{\|V'\|_\infty}{(1-\gamma)^2\gamma^2}} 1,
\end{aligned} \tag{134}$$

where (i) arises from following the routine of (113), (ii) holds by repeating the argument of (114), (iii) follows by taking  $N \geq \frac{\log(\frac{54SAN^2}{(1-\gamma)^\delta})}{(1-\gamma)^2}$  and  $\varepsilon_{\text{opt}} \leq \frac{1-\gamma}{\gamma}$ , and the last inequality holds by  $\|V'\|_\infty \leq \|V^{*,\sigma}\|_\infty \leq \frac{1}{1-\gamma}$ .

Finally, applying Lemma 7 with  $P = \hat{P}^0$  and  $\pi = \hat{\pi}$  yields

$$\|V'\|_\infty \leq \max_{s \in \mathcal{S}} \hat{V}^{\hat{\pi}, \sigma}(s) - \min_{s \in \mathcal{S}} \hat{V}^{\hat{\pi}, \sigma}(s) \leq \frac{1}{\gamma \max\{1-\gamma, \sigma\}},$$

which can be inserted into (134) and gives

$$\left(I - \gamma \underline{P}^{\hat{\pi}, \hat{V}}\right)^{-1} \sqrt{\text{Var}_{\underline{P}^{\hat{\pi}, \hat{V}}}(\hat{V}^{\hat{\pi}, \sigma})} \leq 6\sqrt{\frac{1}{\gamma^3(1-\gamma)^2 \max\{1-\gamma, \sigma\}}} 1 \leq 6\sqrt{\frac{1}{(1-\gamma)^3\gamma^2}} 1.$$

## C Proof of the lower bound with TV distance: Theorem 2

To prove Theorem 2, we shall first construct some hard instances and then characterize the sample complexity requirements over these instances. Note that the hard instances for robust MDPs are different from those for standard MDPs, due to the asymmetric structure induced by the robust RL problem formulation to consider the worst-case performance. By constructing a new class of hard instances inspired by the asymmetric structure of the RMDP, we develop a new lower bound in Theorem 2 that is tighter than prior art (Yang et al., 2022).

### C.1 Construction of the hard problem instances

**Construction of two hard MDPs.** Suppose there are two standard MDPs defined as below:

$$\{\mathcal{M}_\phi = (\mathcal{S}, \mathcal{A}, P^\phi, r, \gamma) \mid \phi = \{0, 1\}\}.$$

Here,  $\gamma$  is the discount parameter,  $\mathcal{S} = \{0, 1, \dots, S-1\}$  is the state space. Given any state  $s \in \{2, 3, \dots, S-1\}$ , the corresponding action space are  $\mathcal{A} = \{0, 1, 2, \dots, A-1\}$ . While for states  $s = 0$  or  $s = 1$ , the action space is only  $\mathcal{A}' = \{0, 1\}$ . For any  $\phi \in \{0, 1\}$ , the transition kernel  $P^\phi$  of the constructed MDP  $\mathcal{M}_\phi$  is defined as

$$P^\phi(s' \mid s, a) = \begin{cases} p\mathbb{1}(s' = 1) + (1-p)\mathbb{1}(s' = 0) & \text{if } (s, a) = (0, \phi) \\ q\mathbb{1}(s' = 1) + (1-q)\mathbb{1}(s' = 0) & \text{if } (s, a) = (0, 1-\phi) \\ \mathbb{1}(s' = 1) & \text{if } s \geq 1 \end{cases}, \tag{135}$$

where  $p$  and  $q$  are set to satisfy

$$0 \leq p \leq 1 \quad \text{and} \quad 0 \leq q = p - \Delta \tag{136}$$

for some  $p$  and  $\Delta > 0$  that shall be introduced later. The above transition kernel  $P^\phi$  implies that state 1 is an absorbing state, namely, the MDP will always stay after it arrives at 1.

Then, we define the reward function as

$$r(s, a) = \begin{cases} 1 & \text{if } s = 1 \\ 0 & \text{otherwise} \end{cases}. \tag{137}$$

Additionally, we choose the following initial state distribution:

$$\varphi(s) = \begin{cases} 1, & \text{if } s = 0 \\ 0, & \text{otherwise} \end{cases}. \quad (138)$$

Here, the constructed two instances are set with different probability transition from state 0 with reward 0 but not state 1 with reward 1 (which were used in standard MDPs (Li et al., 2022b)), yielding a larger gap between the value functions of the two instances.

**Uncertainty set of the transition kernels.** Recalling the uncertainty set assumed throughout this section is defined as  $\mathcal{U}^\sigma(P^\phi)$  with TV distance:

$$\mathcal{U}^\sigma(P) := \mathcal{U}_{\text{TV}}^\sigma(P) = \otimes \mathcal{U}_{\text{TV}}^\sigma(P_{s,a}), \quad \mathcal{U}_{\text{TV}}^\sigma(P_{s,a}) := \left\{ P'_{s,a} \in \Delta(\mathcal{S}) : \frac{1}{2} \|P'_{s,a} - P_{s,a}\|_1 \leq \sigma \right\}, \quad (139)$$

where  $P_{s,a}^\phi := P^\phi(\cdot | s, a)$  is defined similar to (4). In addition, without loss of generality, we recall the radius  $\sigma \in (0, 1 - c_0]$  with  $0 < c_0 < 1$ . With the uncertainty level in hand, taking  $c_1 := \frac{c_0}{2}$ ,  $p$  and  $\Delta$  which determines the instances obey

$$p = (1 + c_1) \max\{1 - \gamma, \sigma\} \quad \text{and} \quad \Delta \leq c_1 \max\{1 - \gamma, \sigma\}, \quad (140)$$

which ensure  $0 \leq p \leq 1$  as follows:

$$(1 + c_1) \sigma \leq 1 - c_0 + c_1 \sigma \leq 1 - \frac{c_0}{2} < 1, \quad (1 + c_1) (1 - \gamma) \leq \frac{3}{2} (1 - \gamma) \leq \frac{3}{4} < 1. \quad (141)$$

Consequently, applying (136) directly leads to

$$p \geq q \geq \max\{1 - \gamma, \sigma\}. \quad (142)$$

To continue, for any  $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$ , we denote the infimum probability of moving to the next state  $s'$  associated with any perturbed transition kernel  $P_{s,a} \in \mathcal{U}^\sigma(P_{s,a}^\phi)$  as

$$\underline{P}^\phi(s' | s, a) := \inf_{P_{s,a} \in \mathcal{U}^\sigma(P_{s,a}^\phi)} P(s' | s, a) = \max\{P(s' | s, a) - \sigma, 0\}, \quad (143)$$

where the last equation can be easily verified by the definition of  $\mathcal{U}^\sigma(P^\phi)$  in (139). As shall be seen, the transition from state 0 to state 1 plays an important role in the analysis, for convenience, we denote

$$\underline{p} := \underline{P}^\phi(1 | 0, \phi) = p - \sigma, \quad \underline{q} := \underline{P}^\phi(1 | 0, 1 - \phi) = q - \sigma, \quad (144)$$

which follows from the fact that  $p \geq q \geq \sigma$  in (142).

**Robust value functions and robust optimal policies.** To proceed, we are ready to derive the corresponding robust value functions, identify the optimal policies, and characterize the optimal values. For any MDP  $\mathcal{M}_\phi$  with the above uncertainty set, we denote  $\pi_\phi^*$  as the optimal policy, and the robust value function of any policy  $\pi$  (resp. the optimal policy  $\pi_\phi^*$ ) as  $V_\phi^{\pi, \sigma}$  (resp.  $V_\phi^{\star, \sigma}$ ). Then, we introduce the following lemma which describes some important properties of the robust (optimal) value functions and optimal policies. The proof is postponed to Appendix C.3.1.

**Lemma 14.** *For any  $\phi = \{0, 1\}$  and any policy  $\pi$ , the robust value function obeys*

$$V_\phi^{\pi, \sigma}(0) = \frac{\gamma(z_\phi^\pi - \sigma)}{(1 - \gamma) \left( 1 + \frac{\gamma(z_\phi^\pi - \sigma)}{1 - \gamma(1 - \sigma)} \right) (1 - \gamma(1 - \sigma))}, \quad (145)$$

where  $z_\phi^\pi$  is defined as

$$z_\phi^\pi := p\pi(\phi | 0) + q\pi(1 - \phi | 0). \quad (146)$$

In addition, the robust optimal value functions and the robust optimal policies satisfy

$$V_{\phi}^{\star,\sigma}(0) = \frac{\gamma(p-\sigma)}{(1-\gamma)\left(1 + \frac{\gamma(p-\sigma)}{1-\gamma(1-\sigma)}\right)(1-\gamma(1-\sigma))}, \quad (147a)$$

$$\pi_{\phi}^{\star}(\phi|s) = 1, \quad \text{for } s \in \mathcal{S}. \quad (147b)$$

## C.2 Establishing the minimax lower bound

Note that our goal is to control the quantity w.r.t. any policy estimator  $\hat{\pi}$  based on the chosen initial distribution  $\varphi$  in (138) and the dataset consisting of  $N$  samples over each state-action pair generated from the nominal transition kernel  $P^{\phi}$ , which gives

$$\langle \varphi, V_{\phi}^{\star,\sigma} - V_{\phi}^{\hat{\pi},\sigma} \rangle = V_{\phi}^{\star,\sigma}(0) - V_{\phi}^{\hat{\pi},\sigma}(0).$$

**Step 1: converting the goal to estimate  $\phi$ .** We make the following useful claim which shall be verified in Appendix C.3.2: With  $\varepsilon \leq \frac{c_1}{32(1-\gamma)}$ , letting

$$\Delta = 32(1-\gamma) \max\{1-\gamma, \sigma\} \varepsilon \leq c_1 \max\{1-\gamma, \sigma\} \quad (148)$$

which satisfies (140), it leads to that for any policy  $\hat{\pi}$ ,

$$\langle \varphi, V_{\phi}^{\star,\sigma} - V_{\phi}^{\hat{\pi},\sigma} \rangle \geq 2\varepsilon(1 - \hat{\pi}(\phi|0)). \quad (149)$$

With this connection established between the policy  $\hat{\pi}$  and its sub-optimality gap as depicted in (149), we can now proceed to build an estimate for  $\phi$ . Here, we denote  $\mathbb{P}_{\phi}$  as the probability distribution when the MDP is  $\mathcal{M}_{\phi}$ , where  $\phi$  can take on values in the set  $\{0, 1\}$ .

Let's assume momentarily that an estimated policy  $\hat{\pi}$  achieves

$$\mathbb{P}_{\phi}\{\langle \varphi, V_{\phi}^{\star,\sigma} - V_{\phi}^{\hat{\pi},\sigma} \rangle \leq \varepsilon\} \geq \frac{7}{8}, \quad (150)$$

then in view of (149), we necessarily have  $\hat{\pi}(\phi|0) \geq \frac{1}{2}$  with probability at least  $\frac{7}{8}$ . With this in mind, we are motivated to construct the following estimate  $\hat{\phi}$  for  $\phi \in \{0, 1\}$ :

$$\hat{\phi} = \arg \max_{a \in \{0,1\}} \hat{\pi}(a|0), \quad (151)$$

which obeys

$$\mathbb{P}_{\phi}\{\hat{\phi} = \phi\} \geq \mathbb{P}_{\phi}\{\hat{\pi}(\phi|0) > 1/2\} \geq \frac{7}{8}. \quad (152)$$

Subsequently, our aim is to demonstrate that (152) cannot occur without an adequate number of samples, which would in turn contradict (149).

**Step 2: probability of error in testing two hypotheses.** Equipped with the aforementioned groundwork, we can now delve into differentiating between the two hypotheses  $\phi \in \{0, 1\}$ . To achieve this, we consider the concept of minimax probability of error, defined as follows:

$$p_e := \inf_{\psi} \max\{\mathbb{P}_0(\psi \neq 0), \mathbb{P}_1(\psi \neq 1)\}. \quad (153)$$

Here, the infimum is taken over all possible tests  $\psi$  constructed from the samples generated from the nominal transition kernel  $P^{\phi}$ .

Moving forward, let us denote  $\mu_{\phi}$  (resp.  $\mu_{\phi}(s)$ ) as the distribution of a sample tuple  $(s_i, a_i, s'_i)$  under the nominal transition kernel  $P^{\phi}$  associated with  $\mathcal{M}_{\phi}$  and the samples are generated independently. Applying



standard results from [Tsybakov \(2009, Theorem 2.2\)](#) and the additivity of the KL divergence (cf. [Tsybakov \(2009, Page 85\)](#)), we obtain

$$\begin{aligned} p_e &\geq \frac{1}{4} \exp \left( -NSAKL(\mu_0 \parallel \mu_1) \right) \\ &= \frac{1}{4} \exp \left\{ -N \left( KL(P^0(\cdot | 0, 0) \parallel P^1(\cdot | 0, 0)) + KL(P^0(\cdot | 0, 1) \parallel P^1(\cdot | 0, 1)) \right) \right\}, \end{aligned} \quad (154)$$

where the last inequality holds by observing that

$$\begin{aligned} KL(\mu_0 \parallel \mu_1) &= \frac{1}{SA} \sum_{s,a,s'} KL(P^0(s' | s, a) \parallel P^1(s' | s, a)) \\ &= \frac{1}{SA} \sum_{a \in \{0,1\}} KL(P^0(\cdot | 0, a) \parallel P^1(\cdot | 0, a)), \end{aligned}$$

Here, the last equality holds by the fact that  $P^0(\cdot | s, a)$  and  $P^1(\cdot | s, a)$  only differ when  $s = 0$ .

Now, our focus shifts towards bounding the terms involving the KL divergence in (154). Given  $p \geq q \geq \max\{1 - \gamma, \sigma\}$  (cf. (142)), applying Lemma 1 (cf. (35)) gives

$$\begin{aligned} KL(P^0(\cdot | 0, 1) \parallel P^1(\cdot | 0, 1)) &= KL(p \parallel q) \leq \frac{(p - q)^2}{(1 - p)p} \stackrel{(i)}{=} \frac{\Delta^2}{p(1 - p)} \\ &\stackrel{(ii)}{=} \frac{1024(1 - \gamma)^2 \max\{1 - \gamma, \sigma\}^2 \varepsilon^2}{p(1 - p)} \\ &\leq \frac{1024(1 - \gamma)^2 \max\{1 - \gamma, \sigma\}^2 \varepsilon^2}{1 - p} \leq \frac{4096}{c_1} (1 - \gamma)^2 \max\{1 - \gamma, \sigma\}^2 \varepsilon^2, \end{aligned} \quad (155)$$

where (i) stems from the definition in (136), (ii) follows by the expression of  $\Delta$  in (148), and the last inequality arises from  $1 - q \geq 1 - p \geq \frac{c_0}{4}$  (see (141)).

Note that it can be shown that  $KL(P^0(\cdot | 0, 0) \parallel P^1(\cdot | 0, 0))$  can be upper bounded in a same manner. Substituting (155) back into (154) demonstrates that: if the sample size is selected as

$$N \leq \frac{c_1 \log 2}{8192(1 - \gamma)^2 \max\{1 - \gamma, \sigma\}^2 \varepsilon^2}, \quad (156)$$

then one necessarily has

$$p_e \geq \frac{1}{4} \exp \left\{ -N \frac{8192}{c_1} (1 - \gamma)^2 \max\{1 - \gamma, \sigma\}^2 \varepsilon^2 \right\} \geq \frac{1}{8}, \quad (157)$$

**Step 3: putting the results together.** Lastly, suppose that there exists an estimator  $\hat{\pi}$  such that

$$\mathbb{P}_0 \{ \langle \varphi, V_0^{*,\sigma} - V_0^{\hat{\pi},\sigma} \rangle > \varepsilon \} < \frac{1}{8} \quad \text{and} \quad \mathbb{P}_1 \{ \langle \varphi, V_1^{*,\sigma} - V_1^{\hat{\pi},\sigma} \rangle > \varepsilon \} < \frac{1}{8}.$$

According to Step 1, the estimator  $\hat{\phi}$  defined in (151) must satisfy

$$\mathbb{P}_0(\hat{\phi} \neq 0) < \frac{1}{8} \quad \text{and} \quad \mathbb{P}_1(\hat{\phi} \neq 1) < \frac{1}{8}.$$

However, this cannot occur under the sample size condition (156) to avoid contradiction with (157). Thus, we have completed the proof.

### C.3 Proof of the auxiliary facts

#### C.3.1 Proof of Lemma 14

**Deriving the robust value function over different states.** For any  $\mathcal{M}_\phi$  with  $\phi \in \{0, 1\}$ , we first characterize the robust value function of any policy  $\pi$  over different states. Before proceeding, we denote the minimum of the robust value function over states as below:

$$V_{\phi, \min}^{\pi, \sigma} := \min_{s \in S} V_\phi^{\pi, \sigma}(s). \quad (158)$$

Clearly, there exists at least one state  $s_{\phi, \min}^\pi$  that satisfies  $V_\phi^{\pi, \sigma}(s_{\phi, \min}^\pi) = V_{\phi, \min}^{\pi, \sigma}$ .

With this in mind, it is easily observed that for any policy  $\pi$ , the robust value function at state  $s = 1$  obeys

$$\begin{aligned} V_\phi^{\pi, \sigma}(1) &= \mathbb{E}_{a \sim \pi(\cdot | 1)} \left[ r(1, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{1, a}^\phi)} \mathcal{P} V_\phi^{\pi, \sigma} \right] \\ &\stackrel{(i)}{=} 1 + \gamma \mathbb{E}_{a \sim \pi(\cdot | 1)} \left[ \underline{P}^\phi(1 | 1, a) V_\phi^{\pi, \sigma}(1) \right] + \gamma \sigma V_{\phi, \min}^{\pi, \sigma} \stackrel{(ii)}{=} 1 + \gamma(1 - \sigma) V_\phi^{\pi, \sigma}(1) + \gamma \sigma V_{\phi, \min}^{\pi, \sigma}, \end{aligned} \quad (159)$$

where (i) holds by  $r(1, a) = 1$  for all  $a \in \mathcal{A}'$  and (143), and (ii) follows from  $P^\phi(1 | 1, a) = 1$  for all  $a \in \mathcal{A}'$ .

Similarly, for any  $s \in \{2, 3, \dots, S-1\}$ , we have

$$\begin{aligned} V_\phi^{\pi, \sigma}(s) &= 0 + \gamma \mathbb{E}_{a \sim \pi(\cdot | s)} \left[ \underline{P}^\phi(1 | s, a) V_\phi^{\pi, \sigma}(1) \right] + \gamma \sigma V_{\phi, \min}^{\pi, \sigma} \\ &= \gamma(1 - \sigma) V_\phi^{\pi, \sigma}(1) + \gamma \sigma V_{\phi, \min}^{\pi, \sigma}, \end{aligned} \quad (160)$$

since  $r(s, a) = 0$  for all  $s \in \{2, 3, \dots, S-1\}$  and the definition in (143).

Finally, we move onto compute  $V_\phi^{\pi, \sigma}(0)$ , the robust value function at state 0 associated with any policy  $\pi$ . First, it obeys

$$\begin{aligned} V_\phi^{\pi, \sigma}(0) &= \mathbb{E}_{a \sim \pi(\cdot | 0)} \left[ r(0, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{0, a}^\phi)} \mathcal{P} V_\phi^{\pi, \sigma} \right] \\ &= 0 + \gamma \pi(\phi | 0) \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{0, \phi}^\phi)} \mathcal{P} V_\phi^{\pi, \sigma} + \gamma \pi(1 - \phi | 0) \inf_{\mathcal{P} \in \mathcal{U}^\sigma(P_{0, 1-\phi}^\phi)} \mathcal{P} V_\phi^{\pi, \sigma}. \end{aligned} \quad (161)$$

Recall the transition kernel defined in (135) and the fact about the uncertainty set over state 0 in (144), it is easily verified that the following probability vector  $P_1 \in \Delta(\mathcal{S})$  obeys  $P_1 \in \mathcal{U}^\sigma(P_{0, \phi}^\phi)$ , which is defined as

$$\begin{aligned} P_1(0) &= 1 - p + \sigma \mathbb{1}(0 = s_{\phi, \min}^\pi), & P_1(1) &= \underline{p} = p - \sigma, \\ P_1(s) &= \sigma \mathbb{1}(s = s_{\phi, \min}^\pi), & \forall s \in \{2, 3, \dots, S-1\}, \end{aligned} \quad (162)$$

where  $\underline{p} = p - \sigma$  due to (144). Similarly, the following probability vector  $P_2 \in \Delta(\mathcal{S})$  also falls into the uncertainty set  $\mathcal{U}^\sigma(P_{0, 1-\phi}^\phi)$ :

$$\begin{aligned} P_2(0) &= 1 - q + \sigma \mathbb{1}(0 = s_{\phi, \min}^\pi), & P_2(1) &= \underline{q} = q - \sigma, \\ P_2(s) &= \sigma \mathbb{1}(s = s_{\phi, \min}^\pi) & \forall s \in \{2, 3, \dots, S-1\}. \end{aligned} \quad (163)$$

It is noticed that  $P_0$  and  $P_1$  defined above are the worst-case perturbations, since the probability mass at state 1 will be moved to the state with the least value. Plugging the above facts about  $P_1 \in \mathcal{U}^\sigma(P_{0, \phi}^\phi)$  and  $P_2 \in \mathcal{U}^\sigma(P_{0, 1-\phi}^\phi)$  into (161), we arrive at

$$\begin{aligned} V_\phi^{\pi, \sigma}(0) &\leq \gamma \pi(\phi | 0) P_1 V_\phi^{\pi, \sigma} + \gamma \pi(1 - \phi | 0) P_2 V_\phi^{\pi, \sigma} \\ &= \gamma \pi(\phi | 0) \left[ (p - \sigma) V_\phi^{\pi, \sigma}(1) + (1 - p) V_\phi^{\pi, \sigma}(0) + \sigma V_{\phi, \min}^{\pi, \sigma} \right] \\ &\quad + \gamma \pi(1 - \phi | 0) \left[ (q - \sigma) V_\phi^{\pi, \sigma}(1) + (1 - q) V_\phi^{\pi, \sigma}(0) + \sigma V_{\phi, \min}^{\pi, \sigma} \right] \\ &\stackrel{(i)}{=} \gamma (z_\phi^\pi - \sigma) V_\phi^{\pi, \sigma}(1) + \gamma \sigma V_{\phi, \min}^{\pi, \sigma} + \gamma (1 - z_\phi^\pi) V_\phi^{\pi, \sigma}(0), \end{aligned} \quad (164)$$

where the last equality holds by the definition of  $z_\phi^\pi$  in (146). To continue, recursively applying (164) yields

$$\begin{aligned} V_\phi^{\pi, \sigma}(0) &\leq \gamma (z_\phi^\pi - \sigma) V_\phi^{\pi, \sigma}(1) + \gamma \sigma V_{\phi, \min}^{\pi, \sigma} + \gamma (1 - z_\phi^\pi) \left[ \gamma (z_\phi^\pi - \sigma) V_\phi^{\pi, \sigma}(1) + \gamma \sigma V_{\phi, \min}^{\pi, \sigma} + \gamma (1 - z_\phi^\pi) V_\phi^{\pi, \sigma}(0) \right] \\ &\stackrel{(i)}{\leq} \gamma (z_\phi^\pi - \sigma) V_\phi^{\pi, \sigma}(1) + \gamma \sigma V_{\phi, \min}^{\pi, \sigma} + \gamma (1 - z_\phi^\pi) \left[ \gamma z_\phi^\pi V_\phi^{\pi, \sigma}(1) + \gamma (1 - z_\phi^\pi) V_\phi^{\pi, \sigma}(0) \right] \end{aligned}$$

$$\begin{aligned}
&\leq \dots \\
&\leq \gamma (z_\phi^\pi - \sigma) V_\phi^{\pi, \sigma}(1) + \gamma \sigma V_{\phi, \min}^{\pi, \sigma} + \gamma z_\phi^\pi \sum_{t=1}^{\infty} \gamma^t (1 - z_\phi^\pi)^t V_\phi^{\pi, \sigma}(1) + \lim_{t \rightarrow \infty} \gamma^t (1 - z_\phi^\pi)^t V_\phi^{\pi, \sigma}(0) \\
&\stackrel{(ii)}{\leq} \gamma (z_\phi^\pi - \sigma) V_\phi^{\pi, \sigma}(1) + \gamma \sigma V_{\phi, \min}^{\pi, \sigma} + \gamma (1 - z_\phi^\pi) \frac{\gamma z_\phi^\pi}{1 - \gamma(1 - z_\phi^\pi)} V_\phi^{\pi, \sigma}(1) + 0 \\
&< \gamma (z_\phi^\pi - \sigma) V_\phi^{\pi, \sigma}(1) + \gamma \sigma V_{\phi, \min}^{\pi, \sigma} + \gamma (1 - z_\phi^\pi) V_\phi^{\pi, \sigma}(1) \\
&= \gamma (1 - \sigma) V_\phi^{\pi, \sigma}(1) + \gamma \sigma V_{\phi, \min}^{\pi, \sigma},
\end{aligned} \tag{165}$$

where (i) uses  $V_{\phi, \min}^{\pi, \sigma} \leq V_\phi^{\pi, \sigma}(1)$ , (ii) follows from  $\gamma(1 - z_\phi^\pi) < 1$ , and the penultimate line follows from the trivial fact that  $\frac{\gamma z_\phi^\pi}{1 - \gamma(1 - z_\phi^\pi)} < 1$ .

Combining (159), (160), and (165), we have that for any policy  $\pi$ ,

$$V_\phi^{\pi, \sigma}(0) = V_{\phi, \min}^{\pi, \sigma}, \tag{166}$$

which directly leads to

$$V_\phi^{\pi, \sigma}(1) = 1 + \gamma (1 - \sigma) V_\phi^{\pi, \sigma}(1) + \gamma \sigma V_{\phi, \min}^{\pi, \sigma} = \frac{1 + \gamma \sigma V_\phi^{\pi, \sigma}(0)}{1 - \gamma (1 - \sigma)}. \tag{167}$$

Let's now return to the characterization of  $V_\phi^{\pi, \sigma}(0)$ . In view of (166), the equality in (164) holds, and we have

$$\begin{aligned}
V_\phi^{\pi, \sigma}(0) &= \gamma (z_\phi^\pi - \sigma) V_\phi^{\pi, \sigma}(1) + \gamma (1 - z_\phi^\pi + \sigma) V_\phi^{\pi, \sigma}(0) \\
&\stackrel{(i)}{=} \gamma (z_\phi^\pi - \sigma) \frac{1 + \gamma \sigma V_\phi^{\pi, \sigma}(0)}{1 - \gamma (1 - \sigma)} + \gamma (1 - z_\phi^\pi + \sigma) V_\phi^{\pi, \sigma}(0) \\
&= \frac{\gamma (z_\phi^\pi - \sigma)}{1 - \gamma (1 - \sigma)} + \gamma \left( 1 + (z_\phi^\pi - \sigma) \frac{\gamma \sigma - (1 - \gamma (1 - \sigma))}{1 - \gamma (1 - \sigma)} \right) V_\phi^{\pi, \sigma}(0) \\
&= \frac{\gamma (z_\phi^\pi - \sigma)}{1 - \gamma (1 - \sigma)} + \gamma \left( 1 - \frac{(1 - \gamma)(z_\phi^\pi - \sigma)}{1 - \gamma (1 - \sigma)} \right) V_\phi^{\pi, \sigma}(0),
\end{aligned}$$

where (i) arises from (167). Solving this relation gives

$$V_\phi^{\pi, \sigma}(0) = \frac{\frac{\gamma (z_\phi^\pi - \sigma)}{1 - \gamma (1 - \sigma)}}{(1 - \gamma) \left( 1 + \frac{\gamma (z_\phi^\pi - \sigma)}{1 - \gamma (1 - \sigma)} \right)}. \tag{168}$$

**The optimal robust policy and optimal robust value function.** We move on to characterize the robust optimal policy and its corresponding robust value function. To begin with, denoting

$$z := \frac{\gamma (z_\phi^\pi - \sigma)}{1 - \gamma (1 - \sigma)}, \tag{169}$$

we rewrite (168) as

$$V_\phi^{\pi, \sigma}(0) = \frac{z}{(1 - \gamma)(1 + z)} =: f(z).$$

Plugging in the fact that  $z_\phi^\pi \geq q \geq \sigma > 0$  in (142), it follows that  $z > 0$ . So for any  $z > 0$ , the derivative of  $f(z)$  w.r.t.  $z$  obeys

$$\frac{(1 - \gamma)(1 + z) - (1 - \gamma)z}{(1 - \gamma)^2(1 + z)^2} = \frac{1}{(1 - \gamma)(1 + z)^2} > 0. \tag{170}$$

Observing that  $f(z)$  is increasing in  $z$ ,  $z$  is increasing in  $z_\phi^\pi$ , and  $z_\phi^\pi$  is also increasing in  $\pi(\phi|0)$  (see the fact  $p \geq q$  in (142)), the optimal policy in state 0 thus obeys

$$\pi_\phi^*(\phi|0) = 1. \quad (171)$$

Considering that the action does not influence the state transition for all states  $s > 0$ , without loss of generality, we choose the robust optimal policy to obey

$$\forall s > 0: \quad \pi_\phi^*(\phi|s) = 1. \quad (172)$$

Taking  $\pi = \pi_\phi^*$ , we complete the proof by showing that the corresponding robust optimal robust value function at state 0 as follows:

$$V_\phi^{*,\sigma}(0) = \frac{\frac{\gamma(z_\phi^{\pi^*} - \sigma)}{1 - \gamma(1 - \sigma)}}{(1 - \gamma) \left(1 + \frac{\gamma(z_\phi^{\pi^*} - \sigma)}{1 - \gamma(1 - \sigma)}\right)} = \frac{\frac{\gamma(p - \sigma)}{1 - \gamma(1 - \sigma)}}{(1 - \gamma) \left(1 + \frac{\gamma(p - \sigma)}{1 - \gamma(1 - \sigma)}\right)}. \quad (173)$$

### C.3.2 Proof of the claim (149)

Plugging in the definition of  $\varphi$ , we arrive at that for any policy  $\pi$ ,

$$\langle \varphi, V_\phi^{*,\sigma} - V_\phi^{\pi,\sigma} \rangle = V_\phi^{*,\sigma}(0) - V_\phi^{\pi,\sigma}(0) = \frac{\frac{\gamma(p - z_\phi^\pi)}{1 - \gamma(1 - \sigma)}}{(1 - \gamma) \left(1 + \frac{\gamma(p - \sigma)}{1 - \gamma(1 - \sigma)}\right) \left(1 + \frac{\gamma(z_\phi^\pi - \sigma)}{1 - \gamma(1 - \sigma)}\right)}, \quad (174)$$

which follows from applying (145) and basic calculus. Then, we proceed to control the above term in two cases separately in terms of the uncertainty level  $\sigma$ .

- When  $\sigma \in (0, 1 - \gamma]$ . Then regarding the important terms in (174), we observe that

$$1 - \gamma < 1 - \gamma(1 - \sigma) \leq 1 - \gamma(1 - (1 - \gamma)) = (1 - \gamma)(1 + \gamma) \leq 2(1 - \gamma), \quad (175)$$

which directly leads to

$$\frac{\gamma(z_\phi^\pi - \sigma)}{1 - \gamma(1 - \sigma)} \stackrel{(i)}{\leq} \frac{\gamma(p - \sigma)}{1 - \gamma(1 - \sigma)} \leq \frac{\gamma c_1(1 - \gamma)}{1 - \gamma(1 - \sigma)} \stackrel{(ii)}{<} c_1 \gamma, \quad (176)$$

where (i) holds by  $z_\phi^\pi < p$ , and (ii) is due to (175). Inserting (175) and (176) back into (174), we arrive at

$$\begin{aligned} \langle \varphi, V_\phi^{*,\sigma} - V_\phi^{\pi,\sigma} \rangle &\geq \frac{\frac{\gamma(p - z_\phi^\pi)}{2(1 - \gamma)}}{(1 - \gamma)(1 + c_1 \gamma)^2} \geq \frac{\gamma(p - z_\phi^\pi)}{8(1 - \gamma)^2} \\ &= \frac{\gamma(p - q)(1 - \pi(\phi|0))}{8(1 - \gamma)^2} = \frac{\gamma \Delta(1 - \pi(\phi|0))}{8(1 - \gamma)^2} \geq 2\varepsilon(1 - \pi(\phi|0)), \end{aligned} \quad (177)$$

where the last inequality holds by setting  $(\gamma \geq 1/2)$

$$\Delta = 32(1 - \gamma)^2 \varepsilon. \quad (178)$$

Finally, it is easily verified that

$$\varepsilon \leq \frac{c_1}{32(1 - \gamma)} \implies \Delta \leq c_1(1 - \gamma).$$

- When  $\sigma \in (1 - \gamma, 1 - c_1]$ . Regarding (174), we observe that

$$\gamma\sigma < 1 - \gamma(1 - \sigma) = 1 - \gamma + \gamma\sigma \leq (1 + \gamma)\sigma \leq 2\sigma, \quad (179)$$

which directly leads to

$$\frac{\gamma(z_\phi^\pi - \sigma)}{1 - \gamma(1 - \sigma)} \leq \frac{\gamma(p - \sigma)}{1 - \gamma(1 - \sigma)} \leq \frac{\gamma c_1 \sigma}{1 - \gamma(1 - \sigma)} \stackrel{(i)}{<} c_1, \quad (180)$$

where (i) holds by (179). Inserting (179) and (180) back into (174), we arrive at

$$\begin{aligned} \langle \varphi, V_\phi^{\star, \sigma} - V_\phi^{\pi, \sigma} \rangle &\geq \frac{\frac{\gamma(p - z_\phi^\pi)}{2\sigma}}{(1 - \gamma)(1 + c_1)^2} \geq \frac{\gamma(p - z_\phi^\pi)}{8(1 - \gamma)\sigma} = \frac{\gamma(p - q)(1 - \pi(\phi|0))}{8(1 - \gamma)\sigma} \\ &= \frac{\gamma\Delta(1 - \pi(\phi|0))}{8(1 - \gamma)\sigma} \geq 2\varepsilon(1 - \pi(\phi|0)), \end{aligned} \quad (181)$$

where the last inequality holds by letting  $(\gamma \geq 1/2)$

$$\Delta = 32(1 - \gamma)\sigma\varepsilon. \quad (182)$$

Finally, it is easily verified that

$$\varepsilon \leq \frac{c_1}{32(1 - \gamma)} \implies \Delta \leq c_1\sigma. \quad (183)$$

## D Proof of the upper bound with $\chi^2$ divergence: Theorem 3

The proof of Theorem 3 mainly follows the structure of the proof of Theorem 1 in Appendix B. Throughout this section, for any nominal transition kernel  $P$ , the uncertainty set is taken as (see (10))

$$\mathcal{U}^\sigma(P) = \mathcal{U}_{\chi^2}^\sigma(P) := \otimes \mathcal{U}_{\chi^2}^\sigma(P_{s,a}), \quad \mathcal{U}_{\chi^2}^\sigma(P_{s,a}) := \left\{ P'_{s,a} \in \Delta(\mathcal{S}) : \sum_{s' \in \mathcal{S}} \frac{(P'(s'|s,a) - P(s'|s,a))^2}{P(s'|s,a)} \leq \sigma \right\}. \quad (184)$$

### D.1 Proof of Theorem 3

In order to control the performance gap  $\|V^{\star, \sigma} - V^{\hat{\pi}, \sigma}\|_\infty$ , recall the error decomposition in (61):

$$V^{\star, \sigma} - V^{\hat{\pi}, \sigma} \leq (V^{\pi^*, \sigma} - \hat{V}^{\pi^*, \sigma}) + \frac{2\gamma\varepsilon_{\text{opt}}}{1 - \gamma} \mathbf{1} + (\hat{V}^{\hat{\pi}, \sigma} - V^{\hat{\pi}, \sigma}), \quad (185)$$

where  $\varepsilon_{\text{opt}}$  (cf. (60)) shall be specified later (which justifies Remark 2). To further control (185), we bound the remaining two terms separately.

**Step 1: controlling  $\|\hat{V}^{\pi^*, \sigma} - V^{\pi^*, \sigma}\|_\infty$ .** Towards this, recall the bound in (66) which holds for any uncertainty set:

$$\begin{aligned} \|\hat{V}^{\pi^*, \sigma} - V^{\pi^*, \sigma}\|_\infty &\leq \gamma \max \left\{ \left\| \left( I - \gamma \hat{\underline{P}}^{\pi^*, \hat{V}} \right)^{-1} \left( \hat{\underline{P}}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right) \right\|_\infty, \right. \\ &\quad \left. \left\| \left( I - \gamma \underline{\hat{P}}^{\pi^*, V} \right)^{-1} \left( \underline{\hat{P}}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right) \right\|_\infty \right\}. \end{aligned} \quad (186)$$

To control the main term  $\hat{\underline{P}}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma}$  in (186), we first introduce an important lemma whose proof is postponed to Appendix D.2.1.

**Lemma 15.** Consider any  $\sigma > 0$  and the uncertainty set  $\mathcal{U}^\sigma(\cdot) := \mathcal{U}_{\chi^2}^\sigma(\cdot)$ . For any  $\delta \in (0, 1)$  and any fixed policy  $\pi$ , one has with probability at least  $1 - \delta$ ,

$$\left\| \hat{\underline{P}}^{\pi, V} V^{\pi, \sigma} - \underline{P}^{\pi, V} V^{\pi, \sigma} \right\|_\infty \leq 4 \sqrt{\frac{2(1 + \sigma) \log(\frac{24SAN}{\delta})}{(1 - \gamma)^2 N}}.$$

Applying Lemma 15 by taking  $\pi = \pi^*$  gives

$$\left\| \hat{\underline{P}}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right\|_\infty \leq 4 \sqrt{\frac{2(1 + \sigma) \log(\frac{24SAN}{\delta})}{(1 - \gamma)^2 N}}, \quad (187)$$

which directly leads to

$$\begin{aligned} & \left\| \left( I - \gamma \hat{\underline{P}}^{\pi^*, \hat{V}} \right)^{-1} \left( \hat{\underline{P}}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right) \right\|_\infty \\ & \leq \left\| \hat{\underline{P}}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right\|_\infty \cdot \left\| \left( I - \gamma \hat{\underline{P}}^{\pi^*, \hat{V}} \right)^{-1} 1 \right\|_\infty \leq 4 \sqrt{\frac{2(1 + \sigma) \log(\frac{24SAN}{\delta})}{(1 - \gamma)^4 N}}. \end{aligned} \quad (188)$$

Similarly, we have

$$\left\| \left( I - \gamma \hat{\underline{P}}^{\pi^*, V} \right)^{-1} \left( \hat{\underline{P}}^{\pi^*, V} V^{\pi^*, \sigma} - \underline{P}^{\pi^*, V} V^{\pi^*, \sigma} \right) \right\|_\infty \leq 4 \sqrt{\frac{2(1 + \sigma) \log(\frac{24SAN}{\delta})}{(1 - \gamma)^4 N}}. \quad (189)$$

Inserting (188) and (189) back to (186) yields

$$\left\| \hat{V}^{\pi^*, \sigma} - V^{\pi^*, \sigma} \right\|_\infty \leq 4 \sqrt{\frac{2(1 + \sigma) \log(\frac{24SAN}{\delta})}{(1 - \gamma)^4 N}}. \quad (190)$$

**Step 2: controlling**  $\left\| \hat{V}^{\hat{\pi}, \sigma} - V^{\hat{\pi}, \sigma} \right\|_\infty$ . Recall the bound in (67) which holds for any uncertainty set:

$$\begin{aligned} \left\| \hat{V}^{\hat{\pi}, \sigma} - V^{\hat{\pi}, \sigma} \right\|_\infty & \leq \gamma \max \left\{ \left\| \left( I - \gamma \hat{\underline{P}}^{\hat{\pi}, V} \right)^{-1} \left( \hat{\underline{P}}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} \right) \right\|_\infty, \right. \\ & \quad \left. \left\| \left( I - \gamma \hat{\underline{P}}^{\hat{\pi}, \hat{V}} \right)^{-1} \left( \hat{\underline{P}}^{\hat{\pi}, V} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{\pi}, V} \hat{V}^{\hat{\pi}, \sigma} \right) \right\|_\infty \right\}. \end{aligned} \quad (191)$$

We introduce the following lemma which controls  $\hat{\underline{P}}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma}$  in (191); the proof is deferred to Appendix D.2.2.

**Lemma 16.** Consider the uncertainty set  $\mathcal{U}^\sigma(\cdot) := \mathcal{U}_{\chi^2}^\sigma(\cdot)$  and any  $\delta \in (0, 1)$ . With probability at least  $1 - \delta$ , one has

$$\left\| \hat{\underline{P}}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - \underline{P}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} \right\|_\infty \leq 12 \sqrt{\frac{2(1 + \sigma) \log(\frac{36SAN^2}{\delta})}{(1 - \gamma)^2 N}} + \frac{2\gamma\epsilon_{\text{opt}}}{1 - \gamma} + 4 \sqrt{\frac{\sigma\epsilon_{\text{opt}}}{(1 - \gamma)^2}}. \quad (192)$$

Repeating the arguments from (187) to (190) yields

$$\left\| \hat{V}^{\hat{\pi}, \sigma} - V^{\hat{\pi}, \sigma} \right\|_\infty \leq 12 \sqrt{\frac{2(1 + \sigma) \log(\frac{36SAN^2}{\delta})}{(1 - \gamma)^4 N}} + \frac{2\gamma\epsilon_{\text{opt}}}{(1 - \gamma)^2} + 4 \sqrt{\frac{\sigma\epsilon_{\text{opt}}}{(1 - \gamma)^4}}. \quad (193)$$

Finally, inserting (190) and (193) back to (185) complete the proof

$$\left\| V^{\pi^*, \sigma} - V^{\hat{\pi}, \sigma} \right\|_\infty \leq \left\| V^{\pi^*, \sigma} - \hat{V}^{\pi^*, \sigma} \right\|_\infty + \frac{2\gamma\epsilon_{\text{opt}}}{1 - \gamma} + \left\| \hat{V}^{\hat{\pi}, \sigma} - V^{\hat{\pi}, \sigma} \right\|_\infty$$

$$\begin{aligned}
&\leq 4\sqrt{\frac{2(1+\sigma)\log(\frac{24SAN}{\delta})}{(1-\gamma)^4N}} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} + 12\sqrt{\frac{2(1+\sigma)\log(\frac{36SAN^2}{\delta})}{(1-\gamma)^4N}} + \frac{2\gamma\varepsilon_{\text{opt}}}{(1-\gamma)^2} + 4\sqrt{\frac{\sigma\varepsilon_{\text{opt}}}{(1-\gamma)^4}} \\
&\leq 24\sqrt{\frac{2(1+\sigma)\log(\frac{36SAN^2}{\delta})}{(1-\gamma)^4N}}, \tag{194}
\end{aligned}$$

where the last line holds by taking  $\varepsilon_{\text{opt}} \leq \min \left\{ \sqrt{\frac{32(1+\sigma)\log(\frac{36SAN^2}{\delta})}{N}}, \frac{4\log(\frac{36SAN^2}{\delta})}{N} \right\}$ .

## D.2 Proof of the auxiliary lemmas

### D.2.1 Proof of Lemma 15

**Step 1: controlling the point-wise concentration.** Consider any fixed policy  $\pi$  and the corresponding robust value vector  $V := V^{\pi, \sigma}$  (independent from  $\hat{P}^0$ ). Invoking Lemma 5 leads to that for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,

$$\begin{aligned}
\left| \hat{P}_{s,a}^{\pi, V} V^{\pi, \sigma} - P_{s,a}^{\pi, V} V^{\pi, \sigma} \right| &= \left| \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left\{ P_{s,a}^0 [V]_{\alpha} - \sqrt{\sigma \text{Var}_{P_{s,a}^0}([V]_{\alpha})} \right\} \right. \\
&\quad \left. - \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left\{ \hat{P}_{s,a}^0 [V]_{\alpha} - \sqrt{\sigma \text{Var}_{\hat{P}_{s,a}^0}([V]_{\alpha})} \right\} \right| \\
&\leq \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left| \left( P_{s,a}^0 - \hat{P}_{s,a}^0 \right) [V]_{\alpha} + \sqrt{\sigma \text{Var}_{\hat{P}_{s,a}^0}([V]_{\alpha})} - \sqrt{\sigma \text{Var}_{P_{s,a}^0}([V]_{\alpha})} \right| \\
&\leq \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left| \left( P_{s,a}^0 - \hat{P}_{s,a}^0 \right) [V]_{\alpha} \right| + \\
&\quad + \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \sqrt{\sigma} \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0}([V]_{\alpha})} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha})} \right|, \tag{195}
\end{aligned}$$

where the first inequality follows by that the maximum operator is 1-Lipschitz, and the second inequality follows from the triangle inequality. Observing that the first term in (195) is exactly the same as (101), recalling the fact in (106) directly leads to: with probability at least  $1 - \delta$ ,

$$\max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left| \left( P_{s,a}^0 - \hat{P}_{s,a}^0 \right) [V]_{\alpha} \right| \leq 2\sqrt{\frac{\log(\frac{2SAN}{\delta})}{(1-\gamma)^2N}} \tag{196}$$

holds for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ . Then the remainder of the proof focuses on controlling the second term in (195).

**Step 2: controlling the second term in (195).** For any given  $(s, a) \in \mathcal{S} \times \mathcal{A}$  and fixed  $\alpha \in [0, \frac{1}{1-\gamma}]$ , applying the concentration inequality (Panaganti and Kalathil, 2022, Lemma 6) with  $\|[V]_{\alpha}\|_{\infty} \leq \frac{1}{1-\gamma}$ , we arrive at

$$\left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0}([V]_{\alpha})} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha})} \right| \leq \sqrt{\frac{2\log(\frac{2}{\delta})}{(1-\gamma)^2N}} \tag{197}$$

holds with probability at least  $1 - \delta$ . To obtain a uniform bound, we first observe the follow lemma proven in Appendix D.2.3.

**Lemma 17.** For any  $V$  obeying  $\|V\|_{\infty} \leq \frac{1}{1-\gamma}$ , the function  $J_{s,a}(\alpha, V) := \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0}([V]_{\alpha})} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha})} \right|$  w.r.t.  $\alpha$  obeys

$$|J_{s,a}(\alpha_1, V) - J_{s,a}(\alpha_2, V)| \leq 4\sqrt{\frac{|\alpha_1 - \alpha_2|}{1-\gamma}}.$$

In addition, we can construct an  $\varepsilon_3$ -net  $N_{\varepsilon_3}$  over  $[0, \frac{1}{1-\gamma}]$  whose size is  $|N_{\varepsilon_3}| \leq \frac{3}{\varepsilon_3(1-\gamma)}$  (Vershynin, 2018). Armed with the above, we can derive the uniform bound over  $\alpha \in [\min_s V(s), \max_s V(s)] \subset [0, 1/(1-\gamma)]$ : with probability at least  $1 - \frac{\delta}{SA}$ , it holds that for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,

$$\begin{aligned}
& \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0}([V]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_\alpha)} \right| \\
& \leq \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0}([V]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_\alpha)} \right| \\
& \stackrel{(i)}{\leq} 4\sqrt{\frac{\varepsilon_3}{1-\gamma}} + \sup_{\alpha \in N_{\varepsilon_3}} \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0}([V]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_\alpha)} \right| \\
& \stackrel{(ii)}{\leq} 4\sqrt{\frac{\varepsilon_3}{1-\gamma}} + \sqrt{\frac{2 \log(\frac{2SA|N_{\varepsilon_3}|}{\delta})}{(1-\gamma)^2 N}} \\
& \stackrel{(iii)}{\leq} 2\sqrt{\frac{2 \log(\frac{2SA|N_{\varepsilon_3}|}{\delta})}{(1-\gamma)^2 N}} \leq 2\sqrt{\frac{2 \log(\frac{24SAN}{\delta})}{(1-\gamma)^2 N}}, \tag{198}
\end{aligned}$$

where (i) holds by the property of  $N_{\varepsilon_3}$ , (ii) follows from (197), (iii) arises from taking  $\varepsilon_3 = \frac{\log(\frac{2SA|N_{\varepsilon_3}|}{\delta})}{8N(1-\gamma)}$ , and the last inequality is verified by  $|N_{\varepsilon_3}| \leq \frac{3}{\varepsilon_3(1-\gamma)} \leq 24N$ .

Inserting (196) and (198) back to (195) and taking the union bound over  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , we arrive at that for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , with probability at least  $1 - \delta$ ,

$$\begin{aligned}
\left| \hat{P}_{s,a}^{\pi,V} V - P_{s,a}^{\pi,V} V \right| & \leq \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left| \left( P_{s,a}^0 - \hat{P}_{s,a}^0 \right) [V]_\alpha \right| + \\
& \quad + \max_{\alpha \in [\min_s V(s), \max_s V(s)]} \left| \sqrt{\sigma \text{Var}_{\hat{P}_{s,a}^0}([V]_\alpha)} - \sqrt{\sigma \text{Var}_{P_{s,a}^0}([V]_\alpha)} \right| \\
& \leq \sqrt{\frac{2 \log(\frac{24SAN}{\delta})}{(1-\gamma)^2 N}} + 2\sqrt{\frac{2\sigma \log(\frac{24SAN}{\delta})}{(1-\gamma)^2 N}} \leq 4\sqrt{\frac{2(1+\sigma) \log(\frac{24SAN}{\delta})}{(1-\gamma)^2 N}}.
\end{aligned}$$

Finally, we complete the proof by recalling the matrix form as below:

$$\left\| \hat{\underline{P}}^{\pi,V} V^{\pi,\sigma} - \underline{P}^{\pi,V} V^{\pi,\sigma} \right\|_\infty \leq \max_{(s,a) \in \mathcal{S} \times \mathcal{A}} \left| \hat{P}_{s,a}^{\pi,V} V - P_{s,a}^{\pi,V} V \right| \leq 4\sqrt{\frac{2(1+\sigma) \log(\frac{24SAN}{\delta})}{(1-\gamma)^2 N}}.$$

### D.2.2 Proof of Lemma 16

**Step 1: decomposing the term of interest.** The proof follows the routine of the proof of Lemma 12 in Appendix B.3.5. To begin with, for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , following the same arguments of (195) yields

$$\begin{aligned}
\left| \hat{\hat{P}}_{s,a}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} - P_{s,a}^{\hat{\pi}, \hat{V}} \hat{V}^{\hat{\pi}, \sigma} \right| & \leq \max_{\alpha \in [\min_s \hat{V}^{\hat{\pi}, \sigma}(s), \max_s \hat{V}^{\hat{\pi}, \sigma}(s)]} \left| \left( P_{s,a}^0 - \hat{P}_{s,a}^0 \right) [\hat{V}^{\hat{\pi}, \sigma}]_\alpha \right| + \\
& \quad + \max_{\alpha \in [\min_s \hat{V}^{\hat{\pi}, \sigma}(s), \max_s \hat{V}^{\hat{\pi}, \sigma}(s)]} \sqrt{\sigma} \left| \sqrt{\text{Var}_{\hat{P}_{s,a}^0}([\hat{V}^{\hat{\pi}, \sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\hat{V}^{\hat{\pi}, \sigma}]_\alpha)} \right|. \tag{199}
\end{aligned}$$

Invoking the fact in (130) (for proving Lemma 12), the first term in (199) obeys

$$\begin{aligned}
\max_{\alpha \in [\min_s \hat{V}^{\hat{\pi}, \sigma}(s), \max_s \hat{V}^{\hat{\pi}, \sigma}(s)]} \left| \left( P_{s,a}^0 - \hat{P}_{s,a}^0 \right) [\hat{V}^{\hat{\pi}, \sigma}]_\alpha \right| & \leq \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \left( P_{s,a}^0 - \hat{P}_{s,a}^0 \right) [\hat{V}^{\hat{\pi}, \sigma}]_\alpha \right| \\
& \leq 4\sqrt{\frac{\log(\frac{3SAN^{3/2}}{(1-\gamma)\delta})}{(1-\gamma)^2 N}} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma}. \tag{200}
\end{aligned}$$

The remainder of the proof will focus on controlling the second term of (199).



**Step 2: controlling the second term of (199).** Towards this, we recall the auxiliary robust MDP  $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$  defined in Appendix B.3.5. Taking the uncertainty set  $\mathcal{U}^\sigma(\cdot) := \mathcal{U}_{\chi^2}^\sigma(\cdot)$  for both  $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$  and  $\widehat{\mathcal{M}}_{\text{rob}}$ , we recall the corresponding robust Bellman operator  $\widehat{\mathcal{T}}_{s,u}^\sigma(\cdot)$  in (119) and the following definition in (120)

$$u^* := \widehat{V}^{*,\sigma}(s) - \gamma \inf_{\mathcal{P} \in \mathcal{U}^\sigma(e_s)} \mathcal{P} \widehat{V}^{*,\sigma}. \quad (201)$$

Following the arguments in Appendix B.3.5, it can be verified that there exists a unique fixed point  $\widehat{Q}_{s,u}^{*,\sigma}$  of the operator  $\widehat{\mathcal{T}}_{s,u}^\sigma(\cdot)$ , which satisfies  $0 \leq \widehat{Q}_{s,u}^{*,\sigma} \leq \frac{1}{1-\gamma} 1$ . In addition, the corresponding robust value function coincides with that of the operator  $\widehat{\mathcal{T}}^\sigma(\cdot)$ , i.e.,  $\widehat{V}_{s,u}^{*,\sigma} = \widehat{V}^{*,\sigma}$ .

We recall the  $N_{\varepsilon_2}$ -net over  $[0, \frac{1}{1-\gamma}]$  whose size obeying  $|N_{\varepsilon_2}| \leq \frac{3}{\varepsilon_2(1-\gamma)}$  (Vershynin, 2018). Then for all  $u \in N_{\varepsilon_2}$  and a fixed  $\alpha$ ,  $\widehat{\mathcal{M}}_{\text{rob}}^{s,u}$  is statistically independent from  $\widehat{P}_{s,a}^0$ , which indicates the independence between  $[\widehat{V}_{s,u}^{*,\sigma}]_\alpha$  and  $\widehat{P}_{s,a}^0$ . With this in mind, invoking the fact in (198) and taking the union bound over all  $(s, a) \in \mathcal{S} \times \mathcal{A}$  and  $u \in N_{\varepsilon_2}$  yields that, with probability at least  $1 - \delta$ ,

$$\max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}_{s,u}^{*,\sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}_{s,u}^{*,\sigma}]_\alpha)} \right| \leq 2 \sqrt{\frac{2 \log(\frac{24SAN|N_{\varepsilon_2}|}{\delta})}{(1-\gamma)^2 N}} \quad (202)$$

holds for all  $(s, a, u) \in \mathcal{S} \times \mathcal{A} \times N_{\varepsilon_2}$ .

To continue, we decompose the term of interest in (199) as follows:

$$\begin{aligned} & \max_{\alpha \in [\min_s \widehat{V}^{\pi, \sigma}(s), \max_s \widehat{V}^{\pi, \sigma}(s)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}^{\pi, \sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}^{\pi, \sigma}]_\alpha)} \right| \\ & \leq \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}^{\pi, \sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}^{\pi, \sigma}]_\alpha)} \right| \\ & \stackrel{(i)}{\leq} \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}^{*, \sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}^{*, \sigma}]_\alpha)} \right| \\ & \quad + \max_{\alpha \in [0, 1/(1-\gamma)]} \left[ \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}^{\pi, \sigma}]_\alpha)} - \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}^{*, \sigma}]_\alpha)} \right| + \left| \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}^{\pi, \sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}^{*, \sigma}]_\alpha)} \right| \right] \\ & \stackrel{(ii)}{\leq} \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}^{*, \sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}^{*, \sigma}]_\alpha)} \right| + \max_{\alpha \in [0, 1/(1-\gamma)]} 2 \sqrt{\frac{2}{(1-\gamma)}} \|\widehat{V}^{\pi, \sigma} - \widehat{V}^{*, \sigma}\|_\alpha \\ & \leq \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}^{*, \sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}^{*, \sigma}]_\alpha)} \right| + 4 \sqrt{\frac{\varepsilon_{\text{opt}}}{(1-\gamma)^2}}, \end{aligned} \quad (203)$$

where (i) holds by the triangle inequality, (ii) arises from applying Lemma 2, and the last inequality holds by (60).

Armed with the above facts, invoking the identity  $\widehat{V}^{*, \sigma} = \widehat{V}_{s,u^*}^{*, \sigma}$  leads to that for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , with probability at least  $1 - \delta$ ,

$$\begin{aligned} & \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}^{*, \sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}^{*, \sigma}]_\alpha)} \right| \\ & = \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}_{s,u^*}^{*, \sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}_{s,u^*}^{*, \sigma}]_\alpha)} \right| \\ & \stackrel{(i)}{\leq} \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}_{s,\bar{u}}^{*, \sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}_{s,\bar{u}}^{*, \sigma}]_\alpha)} \right| \\ & \quad + \max_{\alpha \in [0, 1/(1-\gamma)]} \left[ \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}_{s,u^*}^{*, \sigma}]_\alpha)} - \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}_{s,\bar{u}}^{*, \sigma}]_\alpha)} \right| + \left| \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}_{s,u^*}^{*, \sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}_{s,\bar{u}}^{*, \sigma}]_\alpha)} \right| \right] \\ & \stackrel{(ii)}{\leq} \max_{\alpha \in [0, 1/(1-\gamma)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}_{s,\bar{u}}^{*, \sigma}]_\alpha)} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}_{s,\bar{u}}^{*, \sigma}]_\alpha)} \right| + 4 \sqrt{\frac{\varepsilon_2}{(1-\gamma)}} \end{aligned}$$

$$\begin{aligned}
&\stackrel{\text{(iii)}}{\leq} 2\sqrt{\frac{2\log(\frac{24SAN|N_{\varepsilon_2}|}{\delta})}{(1-\gamma)^2N}} + 4\sqrt{\frac{\varepsilon_2}{(1-\gamma)}} \\
&\leq 6\sqrt{\frac{2\log(\frac{36SAN^2|N_{\varepsilon_2}|}{\delta})}{(1-\gamma)^2N}},
\end{aligned} \tag{204}$$

where (i) holds by the triangle inequality, (ii) arises from applying Lemma 2 and the fact  $\|\widehat{V}_{s,\bar{u}}^{\star,\sigma} - \widehat{V}_{s,u^*}^{\star,\sigma}\|_{\infty} \leq \frac{\varepsilon_2}{(1-\gamma)}$  (see (126)), (iii) follows from (202), and the last inequality holds by letting  $\varepsilon_2 = \frac{2\log(\frac{24SAN|N_{\varepsilon_2}|}{\delta})}{(1-\gamma)N}$ , which leads to  $|N_{\varepsilon_2}| \leq \frac{3}{\varepsilon_2(1-\gamma)} \leq \frac{3N}{2}$ .

In summary, inserting (204) back to (203) and (203) leads to with probability at least  $1 - \delta$ ,

$$\begin{aligned}
&\max_{\alpha \in [\min_s \widehat{V}^{\pi,\sigma}(s), \max_s \widehat{V}^{\pi,\sigma}(s)]} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([\widehat{V}^{\pi,\sigma}]_{\alpha})} - \sqrt{\text{Var}_{P_{s,a}^0}([\widehat{V}^{\pi,\sigma}]_{\alpha})} \right| \\
&\leq 6\sqrt{\frac{2\sigma \log(\frac{36SAN^2|N_{\varepsilon_2}|}{\delta})}{(1-\gamma)^2N}} + 4\sqrt{\frac{\sigma\varepsilon_{\text{opt}}}{(1-\gamma)^2}}
\end{aligned} \tag{205}$$

holds for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ .

**Step 4: finishing up.** Inserting (205) and (200) back to (199), we complete the proof: with probability at least  $1 - \delta$ ,

$$\begin{aligned}
\|\widehat{P}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma} - \underline{P}^{\widehat{\pi}, \widehat{V}} \widehat{V}^{\widehat{\pi}, \sigma}\|_{\infty} &\leq 4\sqrt{\frac{\log(\frac{3SAN^{3/2}}{(1-\gamma)\delta})}{(1-\gamma)^2N}} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} + 6\sqrt{\frac{2\sigma \log(\frac{36SAN^2|N_{\varepsilon_2}|}{\delta})}{(1-\gamma)^2N}} + 4\sqrt{\frac{\sigma\varepsilon_{\text{opt}}}{(1-\gamma)^2}} \\
&\leq 12\sqrt{\frac{2(1+\sigma) \log(\frac{36SAN^2}{\delta})}{(1-\gamma)^2N}} + \frac{2\gamma\varepsilon_{\text{opt}}}{1-\gamma} + 4\sqrt{\frac{\sigma\varepsilon_{\text{opt}}}{(1-\gamma)^2}}.
\end{aligned} \tag{206}$$

### D.2.3 Proof of Lemma 17

For any  $0 \leq \alpha_1, \alpha_2 \leq 1/(1-\gamma)$ , one has

$$\begin{aligned}
&|J_{s,a}(\alpha_1, V) - J_{s,a}(\alpha_2, V)| \\
&= \left| \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([V]_{\alpha_1})} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha_1})} \right| - \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([V]_{\alpha_2})} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha_2})} \right| \right| \\
&\stackrel{\text{(i)}}{\leq} \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([V]_{\alpha_1})} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha_1})} - \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([V]_{\alpha_2})} + \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha_2})} \right| \\
&\leq \left| \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([V]_{\alpha_1})} - \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([V]_{\alpha_2})} \right| + \left| \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha_1})} - \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha_2})} \right| \\
&\stackrel{\text{(ii)}}{\leq} \sqrt{\text{Var}_{\widehat{P}_{s,a}^0}([V]_{\alpha_2}) - \text{Var}_{\widehat{P}_{s,a}^0}([V]_{\alpha_1})} + \sqrt{\text{Var}_{P_{s,a}^0}([V]_{\alpha_2}) - \text{Var}_{P_{s,a}^0}([V]_{\alpha_1})} \\
&\stackrel{\text{(iii)}}{\leq} \sqrt{\left| \widehat{P}_{s,a}^0([V]_{\alpha_1}) \circ ([V]_{\alpha_1}) - ([V]_{\alpha_2}) \circ ([V]_{\alpha_2}) \right|} + \left| \widehat{P}_{s,a}^0([V]_{\alpha_1} + [V]_{\alpha_2}) \cdot \widehat{P}_{s,a}^0([V]_{\alpha_1} - [V]_{\alpha_2}) \right| \\
&\quad + \sqrt{\left| P_{s,a}^0([V]_{\alpha_1}) \circ ([V]_{\alpha_1}) - ([V]_{\alpha_2}) \circ ([V]_{\alpha_2}) \right|} + \left| P_{s,a}^0([V]_{\alpha_1} + [V]_{\alpha_2}) \cdot P_{s,a}^0([V]_{\alpha_1} - [V]_{\alpha_2}) \right| \\
&\leq 2\sqrt{2(\alpha_1 + \alpha_2)|\alpha_1 - \alpha_2|} \leq 4\sqrt{\frac{|\alpha_1 - \alpha_2|}{1-\gamma}}.
\end{aligned} \tag{207}$$

where (i) holds by the fact  $||x| - |y|| \leq |x - y|$  for all  $x, y \in \mathbb{R}$ , (ii) follows from the fact that  $\sqrt{x} - \sqrt{y} \leq \sqrt{x - y}$  for any  $x \geq y \geq 0$  and  $\text{Var}_P([V]_{\alpha_2}) \geq \text{Var}_P([V]_{\alpha_1})$  for any transition kernel  $P \in \Delta(\mathcal{S})$ , (iii) holds by the definition of  $\text{Var}_P(\cdot)$  defined in (36), and the last inequality arises from  $0 \leq \alpha_1, \alpha_2 \leq 1/(1-\gamma)$ .

## E Proof of the lower bound with $\chi^2$ divergence: Theorem 4

To prove Theorem 4, we shall first construct some hard instances and then characterize the sample complexity requirements over these instances. The structure of the hard instances are the same as the ones used in the proof of Theorem 2.

### E.1 Construction of the hard problem instances

First, note that we shall use the same MDPs defined in Appendix C.1 as follows

$$\{\mathcal{M}_\phi = (\mathcal{S}, \mathcal{A}, P^\phi, r, \gamma) \mid \phi = \{0, 1\}\}.$$

In particular, we shall keep the structure of the transition kernel in (135), reward function in (137) and initial state distribution in (138), while  $p$  and  $\Delta$  shall be specified differently later.

**Uncertainty set of the transition kernels.** Recalling the uncertainty set associated with  $\chi^2$  divergence in (184), for any uncertainty level  $\sigma$ , the uncertainty set throughout this section is defined as  $\mathcal{U}^\sigma(P^\phi)$ :

$$\mathcal{U}^\sigma(P^\phi) := \otimes \mathcal{U}_{\chi^2}^\sigma(P_{s,a}^\phi), \quad \mathcal{U}_{\chi^2}^\sigma(P_{s,a}^\phi) := \left\{ P_{s,a} \in \Delta(\mathcal{S}) : \sum_{s' \in \mathcal{S}} \frac{(P(s' | s, a) - P^\phi(s' | s, a))^2}{P^\phi(s' | s, a)} \leq \sigma \right\}. \quad (208)$$

Clearly,  $\mathcal{U}^\sigma(P_{s,a}^\phi) = P_{s,a}^\phi$  whenever the state transition is deterministic for  $\chi^2$  divergence. Here,  $q$  and  $\Delta$  (whose choice will be specified later in more detail) which determine the instances are specified as

$$0 \leq q = \begin{cases} 1 - \gamma & \text{if } \sigma \in (0, \frac{1-\gamma}{4}) \\ \frac{\sigma}{1+\sigma} & \text{if } \sigma \in [\frac{1-\gamma}{4}, \infty) \end{cases}, \quad p = q + \Delta, \quad (209)$$

and

$$0 < \Delta \leq \begin{cases} \frac{1}{4}(1 - \gamma) & \text{if } \sigma \in (0, \frac{1-\gamma}{4}) \\ \min \left\{ \frac{1}{4}(1 - \gamma), \frac{1}{2(1+\sigma)} \right\} & \text{if } \sigma \in [\frac{1-\gamma}{4}, \infty) \end{cases}. \quad (210)$$

This directly ensures that

$$p = \Delta + q \leq \max \left\{ \frac{\frac{1}{2} + \sigma}{1 + \sigma}, \frac{5}{4}(1 - \gamma) \right\} \leq 1$$

since  $\gamma \in [\frac{3}{4}, 1)$ .

To continue, for any  $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$ , we denote the infimum probability of moving to the next state  $s'$  associated with any perturbed transition kernel  $P_{s,a} \in \mathcal{U}^\sigma(P_{s,a}^\phi)$  as

$$\underline{P}^\phi(s' | s, a) := \inf_{P_{s,a} \in \mathcal{U}^\sigma(P_{s,a}^\phi)} P(s' | s, a). \quad (211)$$

In addition, we denote the transition from state 0 to state 1 as follows, which plays an important role in the analysis,

$$\underline{p} := \underline{P}^\phi(1 | 0, \phi), \quad \underline{q} := \underline{P}^\phi(1 | 0, 1 - \phi). \quad (212)$$

Before continuing, we introduce some facts about  $\underline{p}$  and  $\underline{q}$  which are summarized as the following lemma; the proof is postponed to Appendix E.3.1.

**Lemma 18.** *Consider any  $\sigma \in (0, \infty)$  and any  $p, q, \Delta$  obeying (209) and (210), the following properties hold*

$$\begin{cases} \frac{1-\gamma}{2} < \underline{q} < 1 - \gamma, & \underline{q} + \frac{3}{4}\Delta \leq \underline{p} \leq \underline{q} + \Delta \leq \frac{5(1-\gamma)}{4} & \text{if } \sigma \in (0, \frac{1-\gamma}{4}), \\ \underline{q} = 0, & \frac{\sigma+1}{2}\Delta \leq \underline{p} \leq (3+\sigma)\Delta & \text{if } \sigma \in [\frac{1-\gamma}{4}, \infty). \end{cases} \quad (213)$$

**Value functions and optimal policies.** Armed with above facts, we are positioned to derive the corresponding robust value functions, the optimal policies, and its corresponding optimal robust value functions. For any RMDP  $\mathcal{M}_\phi$  with the uncertainty set defined in (208), we denote the robust optimal policy as  $\pi_\phi^*$ , the robust value function of any policy  $\pi$  (resp. the optimal policy  $\pi_\phi^*$ ) as  $V_\phi^{\pi, \sigma}$  (resp.  $V_\phi^{\star, \sigma}$ ). The following lemma describes some key properties of the robust (optimal) value functions and optimal policies whose proof is postponed to Appendix E.3.2.

**Lemma 19.** *For any  $\phi = \{0, 1\}$  and any policy  $\pi$ , one has*

$$V_\phi^{\pi, \sigma}(0) = \frac{\gamma z_\phi^\pi}{(1 - \gamma) \left(1 - \gamma(1 - z_\phi^\pi)\right)}, \quad (214)$$

where  $z_\phi^\pi$  is defined as

$$z_\phi^\pi := \underline{p}\pi(\phi | 0) + \underline{q}\pi(1 - \phi | 0). \quad (215)$$

In addition, the optimal value functions and the optimal policies obey

$$V_\phi^{\star, \sigma}(0) = \frac{\gamma \underline{p}}{(1 - \gamma) (1 - \gamma(1 - \underline{p}))}, \quad (216a)$$

$$\pi_\phi^*(\phi | s) = 1, \quad \text{for } s \in \mathcal{S}. \quad (216b)$$

## E.2 Establishing the minimax lower bound

Our goal is to control the performance gap w.r.t. any policy estimator  $\hat{\pi}$  based on the generated dataset and the chosen initial distribution  $\varphi$  in (138), which gives

$$\langle \varphi, V_\phi^{\star, \sigma} - V_\phi^{\hat{\pi}, \sigma} \rangle = V_\phi^{\star, \sigma}(0) - V_\phi^{\hat{\pi}, \sigma}(0). \quad (217)$$

**Step 1: converting the goal to estimate  $\phi$ .** To achieve the goal, we first introduce the following fact which shall be verified in Appendix E.3.3: given

$$\varepsilon \leq \begin{cases} \frac{1}{72(1-\gamma)} & \text{if } \sigma \in (0, \frac{1-\gamma}{4}), \\ \frac{1}{256(1+\sigma)(1-\gamma)} & \text{if } \sigma \in [\frac{1-\gamma}{4}, \frac{1}{3(1-\gamma)}), \\ \frac{3}{32} & \text{if } \sigma > \frac{1}{3(1-\gamma)}. \end{cases} \quad (218)$$

choosing

$$\Delta = \begin{cases} 18(1 - \gamma)^2 \varepsilon & \text{if } \sigma \in (0, \frac{1-\gamma}{4}), \\ 64(1 + \sigma)(1 - \gamma)^2 \varepsilon & \text{if } \sigma \in [\frac{1-\gamma}{4}, \frac{1}{3(1-\gamma)}), \\ \frac{16}{3(1+\sigma)} \varepsilon & \text{if } \sigma > \frac{1}{3(1-\gamma)}. \end{cases} \quad (219)$$

which satisfies the requirement of  $\Delta$  in (209), it holds that for any policy  $\hat{\pi}$ ,

$$\langle \varphi, V_\phi^{\star, \sigma} - V_\phi^{\hat{\pi}, \sigma} \rangle \geq 2\varepsilon(1 - \hat{\pi}(\phi | 0)). \quad (220)$$

**Step 2: arriving at the final results.** To continue, following the same definitions and argument in Appendix C.2, we recall the minimax probability of the error and its property as follows:

$$p_e \geq \frac{1}{4} \exp \left\{ -N \left( \text{KL}(P^0(\cdot | 0, 0) \| P^1(\cdot | 0, 0)) + \text{KL}(P^0(\cdot | 0, 1) \| P^1(\cdot | 0, 1)) \right) \right\}, \quad (221)$$

then we can complete the proof by showing  $p_e \geq \frac{1}{8}$  given the bound for the sample size  $N$ . In the following, we shall control the KL divergence terms in (221) in three different cases.

- Case 1:  $\sigma \in (0, \frac{1-\gamma}{4})$ . In this case, applying  $\gamma \in [\frac{3}{4}, 1)$  yields

$$\begin{aligned} 1 - q &> 1 - p = 1 - q - \Delta > \gamma - \frac{1-\gamma}{4} > \frac{3}{4} - \frac{1}{16} > \frac{1}{2}, \\ p &\geq q = 1 - \gamma. \end{aligned} \quad (222)$$

Armed with the above facts, applying Lemma 1 (cf. (35)) yields

$$\begin{aligned} \text{KL}(P^0(\cdot | 0, 1) \parallel P^1(\cdot | 0, 1)) &= \text{KL}(p \parallel q) \leq \frac{(p-q)^2}{(1-p)p} \stackrel{(i)}{=} \frac{\Delta^2}{p(1-p)} \\ &\stackrel{(ii)}{=} \frac{324(1-\gamma)^4 \varepsilon^2}{p(1-p)} \\ &\stackrel{(iii)}{\leq} 648(1-\gamma)^3 \varepsilon^2, \end{aligned} \quad (223)$$

where (i) follows from the definition in (209), (ii) holds by plugging in the expression of  $\Delta$  in (219), and (iii) arises from (222). The same bound can be established for  $\text{KL}(P_1^0(\cdot | 0, 0) \parallel P_1^1(\cdot | 0, 0))$ . Substituting (223) back into (221) demonstrates that: if the sample size is chosen as

$$N \leq \frac{\log 2}{1296(1-\gamma)^3 \varepsilon^2}, \quad (224)$$

then one necessarily has

$$p_e \geq \frac{1}{4} \exp \left\{ -N \cdot 1296(1-\gamma)^3 \varepsilon^2 \right\} \geq \frac{1}{8}. \quad (225)$$

- Case 2:  $\sigma \in [\frac{1-\gamma}{4}, \frac{1}{3(1-\gamma)})$ . Applying the facts of  $\Delta$  in (210), one has

$$\begin{aligned} 1 - q &> 1 - p = 1 - q - \Delta \geq \frac{1}{1+\sigma} - \frac{1}{2(1+\sigma)} = \frac{1}{2(1+\sigma)}, \\ p &\geq q = \frac{\sigma}{1+\sigma}. \end{aligned} \quad (226)$$

Given (226), applying Lemma 1 (cf. (35)) yields

$$\begin{aligned} \text{KL}(P^0(\cdot | 0, 1) \parallel P^1(\cdot | 0, 1)) &= \text{KL}(p \parallel q) \leq \frac{(p-q)^2}{(1-p)p} \stackrel{(i)}{=} \frac{\Delta^2}{p(1-p)} \\ &\stackrel{(ii)}{=} \frac{4096(1+\sigma)^2(1-\gamma)^4 \varepsilon^2}{p(1-p)} \\ &\stackrel{(iii)}{\leq} \frac{4096(1+\sigma)^2(1-\gamma)^4 \varepsilon^2}{\frac{\sigma}{2(1+\sigma)^2}} \leq \frac{8192(1-\gamma)^4(1+\sigma)^4 \varepsilon^2}{\sigma}, \end{aligned} \quad (227)$$

where (i) follows from the definition in (209), (ii) holds by plugging in the expression of  $\Delta$  in (219), and (iii) arises from (226). The same bound can be established for  $\text{KL}(P_1^0(\cdot | 0, 0) \parallel P_1^1(\cdot | 0, 0))$ .

Substituting (227) back into (154) demonstrates that: if the sample size is chosen as

$$N \leq \frac{\sigma \log 2}{16384(1-\gamma)^4(1+\sigma)^4 \varepsilon^2}, \quad (228)$$

then one necessarily has

$$p_e \geq \frac{1}{4} \exp \left\{ -N \frac{16384(1-\gamma)^4(1+\sigma)^4 \varepsilon^2}{\sigma} \right\} \geq \frac{1}{8}. \quad (229)$$

- Case 3:  $\sigma > \frac{1}{3(1-\gamma)} \geq \frac{1}{3}$ . Regarding this, one gives

$$\begin{aligned} 1 - q > 1 - p = 1 - q - \Delta &\geq \frac{1}{1 + \sigma} - \frac{1}{4(1 + \sigma)} \geq \frac{1}{2(1 + \sigma)}, \\ p \geq q &\geq \frac{1}{4}. \end{aligned} \quad (230)$$

Given  $p \geq q \geq 1/2$  and (230), applying Lemma 1 (cf. (35)) yields

$$\begin{aligned} \text{KL}(P^0(\cdot | 0, 1) \parallel P^1(\cdot | 0, 1)) &= \text{KL}(p \parallel q) \leq \frac{(p - q)^2}{(1 - p)p} \stackrel{(i)}{=} \frac{\Delta^2}{p(1 - p)} \\ &\stackrel{(ii)}{\leq} \frac{\frac{64}{(1 + \sigma)^2} \varepsilon^2}{p(1 - p)} \\ &\stackrel{(iii)}{\leq} \frac{492 \varepsilon^2}{\sigma}, \end{aligned} \quad (231)$$

where (i) follows from the definition in (209), (ii) holds by plugging in the expression of  $\Delta$  in (219), and (iii) arises from (230). The same bound can be established for  $\text{KL}(P_1^0(\cdot | 0, 0) \parallel P_1^1(\cdot | 0, 0))$ . Substituting (231) back into (154) demonstrates that: if the sample size is chosen as

$$N \leq \frac{\sigma \log 2}{984 \varepsilon^2}, \quad (232)$$

then one necessarily has

$$p_e \geq \frac{1}{4} \exp \left\{ -N \frac{984 \varepsilon^2}{\sigma} \right\} \geq \frac{1}{8}. \quad (233)$$

**Step 3: putting things together.** Finally, summing up the results in (224), (228), and (232), combined with the requirement in (218), one has when

$$\varepsilon \leq c_1 \begin{cases} \frac{1}{1-\gamma} & \text{if } \sigma \in (0, \frac{1-\gamma}{4}) \\ \max \left\{ \frac{1}{(1+\sigma)(1-\gamma)}, 1 \right\} & \text{if } \sigma \in [\frac{1-\gamma}{4}, \infty) \end{cases}, \quad (234)$$

taking

$$N \leq c_2 \begin{cases} \frac{1}{(1-\gamma)^3 \varepsilon^2} & \text{if } \sigma \in (0, \frac{1-\gamma}{4}) \\ \frac{\sigma}{\min\{1, (1-\gamma)^4 (1+\sigma)^4\} \varepsilon^2} & \text{if } \sigma \in [\frac{1-\gamma}{4}, \infty) \end{cases} \quad (235)$$

leads to  $p_e \geq \frac{1}{8}$ , for some universal constants  $c_1, c_2 > 0$ .

### E.3 Proof of the auxiliary facts

We begin with some basic facts about the  $\chi^2$  divergence defined in (34) for any two Bernoulli distributions  $\text{Ber}(w)$  and  $\text{Ber}(x)$ , denoted as

$$f(w, x) := \chi^2(x \parallel w) = \frac{(w - x)^2}{w} + \frac{(1 - w - (1 - x))^2}{1 - w} = \frac{(w - x)^2}{w(1 - w)}. \quad (236)$$

For  $x \in [0, w]$ , it is easily verified that the partial derivative w.r.t.  $x$  obeys  $\frac{\partial f(w, x)}{\partial x} = \frac{2(x - w)}{w(1 - w)} < 0$ , implying that

$$\forall x_1 < x_2 \in [0, w], \quad f(w, x_1) > f(w, x_2). \quad (237)$$

In other words, the  $\chi^2$  divergence  $f(w, x)$  increases as  $x$  decreases from  $w$  to 0.

Next, we introduce the following function for any fixed  $\sigma \in (0, \infty)$  and any  $x \in \left[\frac{\sigma}{1+\sigma}, 1\right)$ :

$$f_\sigma(x) := \inf_{\{y: \chi^2(y||x) \leq \sigma, y \in [0, x]\}} y \stackrel{(i)}{=} \max\left\{0, x - \sqrt{\sigma x(1-x)}\right\} = x - \sqrt{\sigma x(1-x)}, \quad (238)$$

where (i) has been verified in Yang et al. (2022, Corollary B.2), and the last equality holds since  $x \geq \frac{\sigma}{1+\sigma}$ . The next lemma summarizes some useful facts about  $f_\sigma(\cdot)$ , which again has been verified in Yang et al. (2022, Lemma B.12 and Corollary B.2).

**Lemma 20.** *Consider any  $\sigma \in (0, \infty)$ . For  $x \in \left[\frac{\sigma}{1+\sigma}, 1\right)$ ,  $f_\sigma(x)$  is convex and differentiable, which obeys*

$$f'_\sigma(x) = 1 + \frac{\sqrt{\sigma}(2x-1)}{2\sqrt{x(1-x)}}.$$

### E.3.1 Proof of Lemma 18

Let us control  $\underline{q}$  and  $\underline{p}$  respectively.

**Step 1: controlling  $\underline{q}$ .** We shall control  $\underline{q}$  in different cases w.r.t. the uncertainty level  $\sigma$ .

- Case 1:  $\sigma \in (0, \frac{1-\gamma}{4})$ . In this case, recall that  $q = 1 - \gamma$  defined in (209), applying (238) with  $x = q$  leads to

$$1 - \gamma = q > \underline{q} = f_\sigma(q) = 1 - \gamma - \sqrt{\sigma \gamma(1-\gamma)} \geq 1 - \gamma - \sqrt{\frac{1-\gamma}{4} \gamma(1-\gamma)} > \frac{1-\gamma}{2}. \quad (239)$$

- Case 2:  $\sigma \in \left[\frac{1-\gamma}{4}, \infty\right)$ . Note that it suffices to treat  $P_{0,1-\phi}^\phi$  as a Bernoulli distribution  $\text{Ber}(q)$  over states 1 and 0, since we do not allow transition to other states. Recalling  $q = \frac{\sigma}{1+\sigma}$  in (209) and noticing the fact that

$$f(q, 0) = \frac{q^2}{q} + \frac{(1 - (1-q))^2}{1-q} = \frac{q}{(1-q)} = \sigma, \quad (240)$$

one has the probability  $\text{Ber}(0)$  falls into the uncertainty set of  $\text{Ber}(q)$  of size  $\sigma$ . As a result, recalling the definition (212) leads to

$$\underline{q} = P^\phi(1 | 0, 1 - \phi) = 0, \quad (241)$$

since  $\underline{q} \geq 0$ .

**Step 2: controlling  $\underline{p}$ .** To characterize the value of  $\underline{p}$ , we also divide into several cases separately.

- Case 1:  $\sigma \in (0, \frac{1-\gamma}{4})$ . In this case, note that  $p > q = 1 - \gamma \geq \frac{\sigma}{1+\sigma}$ . Therefore, applying that  $f_\sigma(\cdot)$  is convex and the form of its derivative in Lemma 20, one has

$$\begin{aligned} \underline{p} &= f_\sigma(p) \geq f_\sigma(q) + f'_\sigma(q)(p - q) \\ &= \underline{q} + \left(1 + \frac{\sqrt{\sigma}(2q-1)}{2\sqrt{q(1-q)}}\right)\Delta \geq \underline{q} + \left(1 - \frac{\sqrt{\frac{1-\gamma}{4}(1-2(1-\gamma))}}{2\sqrt{(1-\gamma)\gamma}}\right)\Delta \geq \underline{q} + \frac{3\Delta}{4}. \end{aligned} \quad (242)$$

Similarly, applying Lemma 20 leads to

$$\begin{aligned} \underline{p} &= f_\sigma(p) \leq f_\sigma(q) + f'_\sigma(p)(p - q) \\ &= \underline{q} + \left(1 - \frac{\sqrt{\sigma}(1-2p)}{2\sqrt{p(1-p)}}\right)\Delta \leq \underline{q} + \Delta, \end{aligned} \quad (243)$$

where the last inequality holds by  $1 - 2p > 0$  due to the fact  $p = q + \Delta \leq \frac{5}{4}(1 - \gamma) \leq \frac{5}{16} < \frac{1}{2}$  (cf. (210) and  $\gamma \in [\frac{3}{4}, 1)$ ). To sum up, given  $\sigma \in (0, \frac{1-\gamma}{4})$ , combined with (239), we arrive at

$$\underline{q} + \frac{3}{4}\Delta \leq \underline{p} \leq \underline{q} + \Delta \leq \frac{5(1-\gamma)}{4}, \quad (244)$$

where the last inequality holds by  $\Delta \leq \frac{1}{4}(1 - \gamma)$  (see (209)).

- Case 2:  $\sigma \in [\frac{1-\gamma}{4}, \infty)$ . We recall that  $p = q + \Delta > q = \frac{\sigma}{1+\sigma}$  in (209). To derive the lower bound for  $\underline{p}$  in (212), similar to (242), one has

$$\begin{aligned} \underline{p} &= f_\sigma(p) \geq f_\sigma(q) + f'_\sigma(q)(p - q) \\ &= \underline{q} + \left(1 + \frac{\sqrt{\sigma}(2q - 1)}{2\sqrt{q(1 - q)}}\right) \Delta \\ &\stackrel{(i)}{=} 0 + \left(1 + \frac{\sqrt{\sigma}\frac{\sigma-1}{1+\sigma}}{2\sqrt{\frac{\sigma}{1+\sigma} \cdot \frac{1}{1+\sigma}}}\right) \Delta = \left(1 + \frac{\sigma - 1}{2}\right) \Delta = \left(\frac{\sigma + 1}{2}\right) \Delta, \end{aligned} \quad (245)$$

where (i) follows from  $q = \frac{\sigma}{1+\sigma}$  and  $\underline{q} = 0$  (see (241)). For the other direction, similar to (243), we have

$$\begin{aligned} \underline{p} &= f_\sigma(p) \leq f_\sigma(q) + f'_\sigma(p)(p - q) = \underline{q} + \left(1 + \frac{\sqrt{\sigma}(2p - 1)}{2\sqrt{p(1 - p)}}\right) \Delta \\ &\stackrel{(i)}{=} \left(1 + \frac{\sqrt{\sigma}(2p - 1)}{2\sqrt{p(1 - p)}}\right) \Delta \stackrel{(ii)}{=} \left(1 + \frac{\sqrt{\sigma}\left(\frac{\sigma-1}{1+\sigma} + 2\Delta\right)}{2\sqrt{\left(\frac{\sigma}{1+\sigma} + \Delta\right)\left(\frac{1}{1+\sigma} - \Delta\right)}}\right) \Delta \\ &\stackrel{(iii)}{\leq} \left(1 + \frac{\sqrt{\sigma}(1 + 2\Delta)}{2\sqrt{\frac{\sigma}{1+\sigma} \cdot \frac{1}{2(1+\sigma)}}}\right) \Delta \stackrel{(iv)}{\leq} \left(1 + (1 + \sigma)\left(1 + \frac{1}{1 + \sigma}\right)\right) \Delta = (3 + \sigma)\Delta, \end{aligned} \quad (246)$$

where (i) holds by  $\underline{q} = 0$  (see (241)), (ii) follows from plugging in  $p = q + \Delta = \frac{\sigma}{1+\sigma} + \Delta$ , and (iii) and (iv) arises from  $\Delta = \min\left\{\frac{1}{4}(1 - \gamma), \frac{1}{2(1+\sigma)}\right\} \leq 1$  in (210). Combining (245) and (246) yields

$$\frac{\sigma + 1}{2}\Delta \leq \underline{p} \leq (3 + \sigma)\Delta. \quad (247)$$

**Step 3: combining all the results.** Finally, summing up the results for both  $\underline{q}$  (in (239) and (241)) and  $\underline{p}$  (in (244) and (247)), we arrive at the advertised bound.

### E.3.2 Proof of Lemma 19

**The robust value function for any policy  $\pi$ .** For any  $\mathcal{M}_\phi$  with  $\phi \in \{0, 1\}$ , we first characterize the robust value function of any policy  $\pi$  over different states.

Towards this, it is easily observed that for any policy  $\pi$ , the robust value functions at state  $s = 1$  or any  $s \in \{2, 3, \dots, S - 1\}$  obey

$$V_\phi^{\pi, \sigma}(1) \stackrel{(i)}{=} 1 + \gamma V_\phi^{\pi, \sigma}(1) = \frac{1}{1 - \gamma} \quad (248a)$$

and

$$\forall s \in \{2, 3, \dots, S\} : \quad V_\phi^{\pi, \sigma}(s) \stackrel{(ii)}{=} 0 + \gamma V_\phi^{\pi, \sigma}(1) = \frac{\gamma}{1 - \gamma}, \quad (248b)$$



where (i) and (ii) is according to the facts that the transitions defined over states  $s \geq 1$  in (135) give only one possible next state 1, leading to a non-random transition in the uncertainty set associated with  $\chi^2$  divergence, and  $r(1, a) = 1$  for all  $a \in \mathcal{A}'$  and  $r(s, a) = 0$  holds all  $(s, a) \in \{2, 3, \dots, S-1\} \times \mathcal{A}$ .

To continue, the robust value function at state 0 with policy  $\pi$  satisfies

$$\begin{aligned} V_{\phi}^{\pi, \sigma}(0) &= \mathbb{E}_{a \sim \pi(\cdot | 0)} \left[ r(0, a) + \gamma \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{0,a}^{\phi})} \mathcal{P} V_{\phi}^{\pi, \sigma} \right] \\ &= 0 + \gamma \pi(\phi | 0) \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{0,\phi}^{\phi})} \mathcal{P} V_{\phi}^{\pi, \sigma} + \gamma \pi(1 - \phi | 0) \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{0,1-\phi}^{\phi})} \mathcal{P} V_{\phi}^{\pi, \sigma} \end{aligned} \quad (249)$$

$$\stackrel{(i)}{\leq} \frac{\gamma}{1 - \gamma}, \quad (250)$$

where (i) holds by that  $\|V_{\phi}^{\pi, \sigma}\|_{\infty} \leq \frac{1}{1-\gamma}$ . Summing up the results in (248b) and (250) leads to

$$\forall s \in \{2, 3, \dots, S\}, \quad V_{\phi}^{\pi, \sigma}(1) > V_{\phi}^{\pi, \sigma}(s) \geq V_{\phi}^{\pi, \sigma}(0). \quad (251)$$

With the transition kernel in (135) over state 0 and the fact in (251), (249) can be rewritten as

$$\begin{aligned} V_{\phi}^{\pi, \sigma}(0) &= \gamma \pi(\phi | 0) \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{0,\phi}^{\phi})} \mathcal{P} V_{\phi}^{\pi, \sigma} + \gamma \pi(1 - \phi | 0) \inf_{\mathcal{P} \in \mathcal{U}^{\sigma}(P_{0,1-\phi}^{\phi})} \mathcal{P} V_{\phi}^{\pi, \sigma} \\ &\stackrel{(i)}{=} \gamma \pi(\phi | 0) \left[ \underline{p} V_{\phi}^{\pi, \sigma}(1) + (1 - \underline{p}) V_{\phi}^{\pi, \sigma}(0) \right] + \gamma \pi(1 - \phi | 0) \left[ \underline{q} V_{\phi}^{\pi, \sigma}(1) + (1 - \underline{q}) V_{\phi}^{\pi, \sigma}(0) \right] \\ &\stackrel{(ii)}{=} \gamma z_{\phi}^{\pi} V_{\phi}^{\pi, \sigma}(1) + \gamma (1 - z_{\phi}^{\pi}) V_{\phi}^{\pi, \sigma}(0) \\ &= \frac{\gamma z_{\phi}^{\pi}}{(1 - \gamma) (1 - \gamma (1 - z_{\phi}^{\pi}))}, \end{aligned} \quad (252)$$

where (i) holds by the definition of  $\underline{p}$  and  $\underline{q}$  in (212), (ii) follows from the definition of  $z_{\phi}^{\pi}$  in (215), and the last line holds by applying (248a) and solving the resulting linear equation for  $V_{\phi}^{\pi, \sigma}(0)$ .

**Optimal policy and its optimal value function.** To continue, observing that  $V_{\phi}^{\pi, \sigma}(0) =: f(z_{\phi}^{\pi})$  is increasing in  $z_{\phi}^{\pi}$  since the derivative of  $f(z_{\phi}^{\pi})$  w.r.t.  $z_{\phi}^{\pi}$  obeys

$$f'(z_{\phi}^{\pi}) = \frac{\gamma(1 - \gamma) (1 - \gamma(1 - z_{\phi}^{\pi})) - \gamma^2 z_{\phi}^{\pi} (1 - \gamma)}{(1 - \gamma)^2 (1 - \gamma(1 - z_{\phi}^{\pi}))^2} = \frac{\gamma}{(1 - \gamma(1 - z_{\phi}^{\pi}))^2} > 0,$$

where the last inequality holds by  $0 \leq z_{\phi}^{\pi} \leq 1$ . Further,  $z_{\phi}^{\pi}$  is also increasing in  $\pi(\phi | 0)$  (see the fact  $\underline{p} \geq \underline{q}$  in (212)), the optimal robust policy in state 0 thus obeys

$$\pi_{\phi}^{\star}(\phi | 0) = 1. \quad (253)$$

Considering that the action does not influence the state transition for all states  $s > 0$ , without loss of generality, we choose the optimal robust policy to obey

$$\forall s > 0: \quad \pi_{\phi}^{\star}(\phi | s) = 1. \quad (254)$$

Taking  $\pi = \pi_{\phi}^{\star}$  and  $z_{\phi}^{\pi_{\phi}^{\star}} = \underline{p}$  in (252), we complete the proof by showing the corresponding optimal robust value function at state 0 as follows:

$$V_{\phi}^{\star, \sigma}(0) = \frac{\gamma z_{\phi}^{\pi_{\phi}^{\star}}}{(1 - \gamma) (1 - \gamma (1 - z_{\phi}^{\pi_{\phi}^{\star}}))} = \frac{\gamma \underline{p}}{(1 - \gamma) (1 - \gamma (1 - \underline{p}))}.$$

### E.3.3 Proof of the claim (220)

Plugging in the definition of  $\varphi$ , we arrive at that for any policy  $\pi$ ,

$$\begin{aligned} \langle \varphi, V_{\phi}^{\star, \sigma} - V_{\phi}^{\pi, \sigma} \rangle &= V_{\phi}^{\star, \sigma}(0) - V_{\phi}^{\pi, \sigma}(0) \\ &\stackrel{(i)}{=} \frac{\gamma \underline{p}}{(1-\gamma)(1-\gamma(1-\underline{p}))} - \frac{\gamma z_{\phi}^{\pi}}{(1-\gamma)(1-\gamma(1-z_{\phi}^{\pi}))} \\ &= \frac{\gamma(\underline{p} - z_{\phi}^{\pi})}{(1-\gamma(1-\underline{p}))(1-\gamma(1-z_{\phi}^{\pi}))} \stackrel{(ii)}{\geq} \frac{\gamma(\underline{p} - z_{\phi}^{\pi})}{(1-\gamma(1-\underline{p}))^2} \stackrel{(iii)}{=} \frac{\gamma(\underline{p} - \underline{q})(1-\pi(\phi|0))}{(1-\gamma(1-\underline{p}))^2}, \end{aligned} \quad (255)$$

where (i) holds by applying Lemma 19, (ii) arises from  $z_{\phi}^{\pi} \leq \underline{p}$  (see the definition of  $z_{\phi}^{\pi}$  in (215) and the fact  $\underline{p} \geq \underline{q} + \frac{3\Delta}{4}$  in (212)), and (iii) follows from the definition of  $z_{\phi}^{\pi}$  in (215).

To further control (255), we consider it in two cases separately:

- Case 1:  $\sigma \in (0, \frac{1-\gamma}{4})$ . In this case, applying Lemma 18 to (255) yields

$$\begin{aligned} \langle \varphi, V_{\phi}^{\star, \sigma} - V_{\phi}^{\pi, \sigma} \rangle &\geq \frac{\gamma(\underline{p} - \underline{q})(1-\pi(\phi|0))}{(1-\gamma(1-\underline{p}))^2} \geq \frac{\gamma \frac{3\Delta}{4}(1-\pi(\phi|0))}{\left(1-\gamma\left(1-\frac{5(1-\gamma)}{4}\right)\right)^2} \\ &\geq \frac{\Delta(1-\pi(\phi|0))}{9(1-\gamma)^2} = 2\varepsilon(1-\pi(\phi|0)), \end{aligned} \quad (256)$$

where the penultimate inequality follows from  $\gamma \geq 3/4$ , and the last inequality holds by taking the specification of  $\Delta$  in (219) as follows:

$$\Delta = 18(1-\gamma)^2\varepsilon. \quad (257)$$

It is easily verified that taking  $\varepsilon \leq \frac{1}{72(1-\gamma)}$  as in (218) directly leads to meeting the requirement in (210), i.e.,  $\Delta \leq \frac{1}{4}(1-\gamma)$ .

- Case 2:  $\sigma \in [\frac{1-\gamma}{4}, \infty)$ . Similarly, applying Lemma 18 to (255) gives

$$\langle \varphi, V_{\phi}^{\star, \sigma} - V_{\phi}^{\pi, \sigma} \rangle \geq \frac{\gamma(\underline{p} - \underline{q})(1-\pi(\phi|0))}{(1-\gamma(1-\underline{p}))^2} \geq \frac{\gamma \frac{\sigma+1}{2} \Delta(1-\pi(\phi|0))}{\min\{1, (1-\gamma(1-(3+\sigma)\Delta))^2\}} \quad (258)$$

Before continuing, it can be verified that

$$\begin{aligned} 1-\gamma(1-(3+\sigma)\Delta) &= 1-\gamma+\gamma(3+\sigma)\Delta \stackrel{(i)}{\leq} 1-\gamma+(3+\sigma)\min\left\{\frac{1}{4}(1-\gamma), \frac{1}{2(\sigma+1)}\right\} \\ &\leq \min\left\{2(1+\sigma)(1-\gamma), \frac{3}{2}\right\}, \end{aligned} \quad (259)$$

where (i) is obtained by  $\Delta \leq \min\left\{\frac{1}{4}(1-\gamma), \frac{1}{2(1+\sigma)}\right\}$  (see (209)). Applying the above fact to (258) gives

$$\begin{aligned} \langle \varphi, V_{\phi}^{\star, \sigma} - V_{\phi}^{\pi, \sigma} \rangle &\geq \frac{\gamma \frac{\sigma+1}{2} \Delta(1-\pi(\phi|0))}{\min\{1, (1-\gamma(1-(3+\sigma)\Delta))^2\}} \stackrel{(i)}{\geq} \frac{3(\sigma+1)\Delta(1-\pi(\phi|0))}{8\min\{4(1+\sigma)^2(1-\gamma)^2, 1\}} \\ &\geq \frac{\Delta(1-\pi(\phi|0))}{\min\{32(1+\sigma)(1-\gamma)^2, \frac{8}{3(1+\sigma)}\}} = 2\varepsilon(1-\pi(\phi|0)), \end{aligned} \quad (260)$$

where (i) holds by  $\gamma \geq \frac{3}{4}$  and (258), and the last equality holds by the specification in (219):

$$\Delta = \begin{cases} 64(1+\sigma)(1-\gamma)^2\varepsilon & \text{if } \sigma \in \left[\frac{1-\gamma}{4}, \frac{1}{3(1-\gamma)}\right), \\ \frac{16}{3(1+\sigma)}\varepsilon & \text{if } \sigma > \frac{1}{3(1-\gamma)}. \end{cases} \quad (261)$$

As a result, it is easily verified that the requirement in (210)

$$\Delta \leq \min \left\{ \frac{1}{4}(1-\gamma), \frac{1}{2(1+\sigma)} \right\} \quad (262)$$

is met if we let

$$\varepsilon \leq \begin{cases} \frac{1}{256(1+\sigma)(1-\gamma)} & \text{if } \sigma \in \left[\frac{1-\gamma}{4}, \frac{1}{3(1-\gamma)}\right), \\ \frac{3}{32} & \text{if } \sigma > \frac{1}{3(1-\gamma)}, \end{cases} \quad (263)$$

as in (218).

The proof is then completed by summing up the results in the above two cases.